

Calcolo di Matrici a Spirale ad Altissime Prestazioni: Dalle Limitazioni Fisiche della Memoria all'Esecuzione Parallela Sub-Secondo su GPU

Advanced Agentic Coding Team

June 5, 2026

Abstract

La generazione e l'analisi di matrici a spirale di grandi dimensioni costituiscono un problema classico dell'informatica algoritmica, caratterizzato da sfide significative in termini di complessità computazionale e gestione delle risorse hardware. Quando le dimensioni della matrice scalano a centinaia di migliaia di righe e colonne, le implementazioni sequenziali convenzionali falliscono a causa di colli di bottiglia insormontabili legati alla memoria (Out-of-Memory). Questo articolo presenta uno studio approfondito sull'ottimizzazione del calcolo di matrici a spirale rettangolari. Dimostriamo la transizione da un approccio ingenuo in Python caratterizzato da una complessità spaziale $O(M \times N)$ a un modello matematico ad indice diretto operante in spazio $O(1)$. Successivamente, esploriamo l'accelerazione di tale modello mediante tecniche di calcolo parallelo su CPU multi-core tramite OpenMP e, infine, mediante programmazione massivamente parallela su GPU NVIDIA tramite l'architettura CUDA. I risultati empirici dimostrano come il calcolo di una matrice a spirale da 30 miliardi di elementi, che originariamente portava al crash del sistema per esaurimento di memoria, possa essere eseguito in un tempo inferiore ad un singolo secondo, evidenziando l'importanza di un design software strettamente allineato alle caratteristiche dell'hardware sottostante.

Contents

1	Introduzione e Contesto Storico	3
1.1	Applicazioni Pratiche delle Traiettorie a Spirale	3
1.2	La Sfida della Complessità	3
2	Il Modello Matematico $O(1)$ in Spazio (Il “Trucco $O(1)$”)	5
2.1	Identificazione del Layer Concentrico	5
2.2	Dimensioni del Sotto-Rettangolo	5
2.3	Determinazione del Valore Iniziale del Layer	5
2.4	Calcolo del Valore di Cella in Base al Bordo	5
2.5	Esempio Numerico Dettagliato	6
3	Limiti Fisici della Memoria e Analisi dell’Architettura CPython	8
3.1	L’Anatomia di un Intero in CPython	8
3.2	Analisi dei Puntatori e Struttura List-of-Lists	8
3.3	Confronto con Tipi Nativi C	9
4	Parallelizzazione CPU con OpenMP	10
4.1	Direttive di Controllo e Collapsing dei Cicli	10
4.2	Riduzione XOR e Prevenzione delle Condizioni di Corsa	10
5	Ottimizzazione Massivamente Parallela su GPU tramite CUDA	12
5.1	Architettura Hardware e Organizzazione dei Thread	12
5.2	Riduzione Intra-Block tramite Warp Shuffle	12
5.3	Esecuzione Branchless e Prevenzione della Divergenza dei Warp	13
6	Codici di Riferimento Completi	15
6.1	Implementazione in Python 3 (Senza Allocazione Heap)	15
6.2	Implementazione in C con Parallelismo OpenMP	16
6.3	Implementazione GPU Massivamente Parallela con CUDA	17
7	Risultati Sperimentali e Benchmarking	20
7.1	Configurazione dell’Ambiente di Test	20
7.2	Analisi dei Tempi di Esecuzione	20
7.3	Analisi dei Risultati	20
7.4	Validazione Matematica del Checksum XOR	21
8	Discussione, Limitazioni e Sviluppi Futuri	23
8.1	L’Equivoco Ontologico del Checksum e Payload Reali	23
8.2	Località della Memoria e Accessi non Coalescenti: Shared Memory Tiling	23
8.3	Estensione Multidimensionale e Sistemi Distribuiti	24

1 Introduzione e Contesto Storico

Una matrice a spirale è una griglia bidimensionale di dimensioni $M \times N$ popolata da valori incrementali sequenziali disposti secondo un percorso continuo che inizia dall'angolo superiore sinistro, si muove verso destra lungo il bordo esterno, e successivamente ruota in senso orario ripiegando verso l'interno ad ogni rivoluzione completa. Dal punto di vista algoritmico elementare, la costruzione di tale struttura viene spesso assegnata come esercizio accademico per verificare la comprensione dei costrutti di controllo di flusso e della manipolazione degli indici di array. Tuttavia, quando le dimensioni del problema crescono fino a contesti industriali o scientifici, l'approccio standard mostra limiti invalicabili legati all'architettura dei calcolatori moderni.

1.1 Applicazioni Pratiche delle Traiettorie a Spirale

Sebbene possa apparire come un problema puramente matematico o ricreativo, il pattern di scansione e generazione a spirale trova importanti riscontri in svariati ambiti tecnologici e scientifici:

- **Imaging a Risonanza Magnetica (MRI) e Tomografia Computerizzata (CT):** Nelle moderne apparecchiature di diagnostica medica, l'acquisizione dei dati nel dominio della frequenza (spazio k) viene frequentemente eseguita seguendo traiettorie a spirale. Questo approccio riduce drasticamente i tempi di scansione rispetto alle griglie cartesiane tradizionali, minimizzando la sensibilità ai movimenti del paziente e ottimizzando l'efficienza dei gradienti magnetici.
- **Memorizzazione Fisica su Supporti Ottici e Magnetici:** I dischi rigidi (HDD) e i supporti ottici come CD, DVD e Blu-ray registrano fisicamente le informazioni lungo una traccia a spirale continua. Comprendere la mappatura geometrica tra coordinate cartesiane della testina e indice lineare della spirale è fondamentale per garantire un accesso rapido ed efficiente ai dati.
- **Curve Space-Filling e Serializzazione di Dati:** Algoritmi di compressione d'immagine ed elaborazione di segnali bidimensionali utilizzano curve a spirale (o varianti come la curva di Hilbert e di Peano) per preservare la coerenza spaziale e la località dei dati durante la linearizzazione di informazioni multidimensionali.
- **Computer Vision e Grafica Computazionale:** Nella rasterizzazione o nel ray-tracing di scene complesse, l'analisi ordinata dal centro dell'immagine verso l'esterno consente di elaborare con priorità assoluta le porzioni focali dell'inquadratura, migliorando l'esperienza utente in applicazioni di rendering interattivo o di compressione video adattiva.

1.2 La Sfida della Complessità

Il metodo tradizionale per generare una matrice a spirale prevede la simulazione fisica della traiettoria di scrittura. Questo processo impiega puntatori di direzione (destra, basso, sinistra, alto) e aggiorna ciclicamente le coordinate correnti all'interno di un ciclo che si ripete per $M \times N$ volte. Sebbene la complessità temporale sia linearmente proporzionale al numero di elementi, ovvero $O(M \times N)$, la necessità di mantenere in memoria la griglia bidimensionale impone una complessità spaziale altrettanto pari a $O(M \times N)$.

Quando si affronta una matrice di dimensioni 100.000×300.000 , pari a 30 miliardi di celle, i requisiti di memoria per l'allocazione della struttura dati superano abbondantemente le capacità fisiche dei sistemi di calcolo convenzionali, causando crash del software o costringendo il sistema operativo a fare affidamento sulla memoria virtuale su disco (paging), azzerando le prestazioni complessive.

In questo articolo analizzeremo in dettaglio come risolvere questo problema rimuovendo totalmente il requisito di memoria $O(M \times N)$ attraverso un approccio matematico ad indice diretto $O(1)$ in spazio. Dimostreremo poi come questo modello, liberato dai vincoli di I/O e allocazione dinamica della memoria, possa essere parallelizzato su CPU multi-core e su GPU massivamente parallele per minimizzare il tempo di calcolo complessivo.

2 Il Modello Matematico $O(1)$ in Spazio (Il “Trucco $O(1)$ ”)

La chiave per superare i limiti di memoria risiede nella possibilità di calcolare il valore di una generica cella in posizione (r, c) all'interno di una matrice $M \times N$ in tempo costante $O(1)$, in modo del tutto indipendente dalla conoscenza dei valori memorizzati nelle celle adiacenti. Questo approccio elimina la necessità di dichiarare ed allocare la matrice in memoria, poiché ogni cella può essere calcolata “al volo” (on-the-fly) quando richiesto.

2.1 Identificazione del Layer Concentrico

Il percorso a spirale può essere modellato come una serie di anelli rettangolari concentrici (detti *layer*). Per una cella situata alla riga r (da 0 a $M - 1$) e alla colonna c (da 0 a $N - 1$), il livello del layer L al quale la cella appartiene è determinato dalla distanza minima della cella stessa dai quattro confini esterni della matrice:

$$L = \min(r, c, M - 1 - r, N - 1 - c) \quad (1)$$

Un valore di $L = 0$ indica che la cella appartiene al bordo più esterno della matrice; all'aumentare di L , la cella si trova in layer progressivamente più interni.

2.2 Dimensioni del Sotto-Rettangolo

Ciascun layer L definisce il perimetro di un sotto-rettangolo interno alla matrice originaria. La larghezza w e l'altezza h di questo sotto-rettangolo sono linearmente dipendenti dal livello del layer L e dalle dimensioni iniziali M e N :

$$w = N - 2L \quad (2)$$

$$h = M - 2L \quad (3)$$

Questi valori rappresentano l'estensione del rettangolo corrente sul quale stiamo calcolando i valori degli elementi di bordo.

2.3 Determinazione del Valore Iniziale del Layer

Ogni layer L ha un valore iniziale S , che corrisponde al valore del suo elemento d'angolo superiore sinistro. Poiché la spirale viene popolata sequenzialmente dall'esterno verso l'interno, il valore S del layer L è pari al numero totale di elementi contenuti in tutti i layer esterni più esterni $(0, 1, \dots, L - 1)$, incrementato di una unità per convenzione di indicizzazione a partire da 1.

Questo conteggio può essere espresso geometricamente con estrema semplicità: la quantità di elementi nei layer esterni è pari all'area totale della matrice iniziale meno l'area del sotto-rettangolo associato al layer L corrente. Otteniamo quindi:

$$S = (M \times N) - (w \times h) + 1 \quad (4)$$

Questa formulazione matematica evita del tutto la necessità di sommare iterativamente i perimetri dei singoli anelli esterni, riducendo l'operazione a una semplice differenza di aree eseguibile in poche operazioni aritmetiche fondamentali.

2.4 Calcolo del Valore di Cella in Base al Bordo

Una volta determinato il valore di partenza del layer S , la larghezza w , l'altezza h , e l'indice del layer L , il valore esatto della cella $V(r, c)$ dipende da quale dei quattro bordi (lati) del sotto-rettangolo ospita la coordinata (r, c) . Analizziamo le quattro partizioni in modo esaustivo:

1. **Bordo Superiore** ($r = L$): La riga coincide con l'indice del layer. La colonna si muove da sinistra a destra nel range $[L, L + w - 1]$. Il valore della cella è pari al valore iniziale S più lo spostamento lineare lungo la colonna:

$$V(r, c) = S + (c - L) \quad (5)$$

2. **Bordo Destro** ($c = L + w - 1$): La colonna coincide con il limite destro del sotto-rettangolo. La riga scende da $L + 1$ a $L + h - 1$. Il valore viene accumulato sommando la larghezza intera del bordo superiore w (già percorsa) e lo spostamento lungo la riga:

$$V(r, c) = S + w + (r - L) - 1 \quad (6)$$

3. **Bordo Inferiore** ($r = L + h - 1$ e $h > 1$): La riga coincide con il limite inferiore del sotto-rettangolo. La colonna si muove da destra verso sinistra nel range $[L, L + w - 2]$. Il valore si calcola sommando gli elementi dei bordi superiore e destro già completati ($w + h - 2$), e lo spostamento all'indietro a partire dal bordo destro:

$$V(r, c) = S + w + h - 2 + (L + w - 1 - c) \quad (7)$$

4. **Bordo Sinistro (Altrimenti)**: La colonna coincide con l'indice del layer L , e la riga sale dal basso verso l'alto nel range $[L + 1, L + h - 2]$. Il valore accumula l'intera lunghezza del bordo superiore, destro e inferiore ($2w + h - 3$), più la distanza percorsa in salita a partire dall'angolo inferiore sinistro:

$$V(r, c) = S + 2w + h - 3 + (L + h - 1 - r) \quad (8)$$

Grazie a questa tassonomia, qualsiasi coordinata (r, c) può essere immediatamente risolta fornendo il valore corretto in tempo di calcolo costante $O(1)$ e senza richiedere allocazioni supplementari di memoria.

2.5 Esempio Numerico Dettagliato

Per validare formalmente il modello matematico, consideriamo una matrice rettangolare di dimensioni $M = 4$ (righe) e $N = 5$ (colonne). La rappresentazione tabellare della matrice a spirale attesa è la seguente:

$$A = \begin{pmatrix} 1 & 2 & 3 & 4 & 5 \\ 14 & 15 & 16 & 17 & 6 \\ 13 & 20 & 19 & 18 & 7 \\ 12 & 11 & 10 & 9 & 8 \end{pmatrix} \quad (9)$$

Applichiamo il nostro modello ad alcune celle campione per verificarne l'esattezza matematica:

- **Caso 1: Cella** $(0, 0)$

- Calcolo del layer: $L = \min(0, 0, 4 - 1 - 0, 5 - 1 - 0) = 0$.
- Dimensioni: $w = 5 - 2(0) = 5$, $h = 4 - 2(0) = 4$.
- Valore iniziale: $S = (4 \times 5) - (5 \times 4) + 1 = 1$.
- Posizionamento: $r = 0 = L$ (Bordo Superiore).
- Valore: $V(0, 0) = S + (c - L) = 1 + (0 - 0) = 1$. **Corretto.**

- **Caso 2: Cella** $(2, 4)$

- Calcolo del layer: $L = \min(2, 4, 4 - 1 - 2, 5 - 1 - 4) = \min(2, 4, 1, 0) = 0$.

- Dimensioni: $w = 5, h = 4$. Valore iniziale: $S = 1$.
- Posizionamento: $c = 4 = L + w - 1$ ($0 + 5 - 1 = 4$) (Bordo Destro).
- Valore: $V(2,4) = S + w + (r - L) - 1 = 1 + 5 + (2 - 0) - 1 = 7$. **Corretto.**

• **Caso 3: Cella (3,1)**

- Calcolo del layer: $L = \min(3, 1, 4 - 1 - 3, 5 - 1 - 1) = \min(3, 1, 0, 3) = 0$.
- Dimensioni: $w = 5, h = 4$. Valore iniziale: $S = 1$.
- Posizionamento: $r = 3 = L + h - 1$ ($0 + 4 - 1 = 3$) (Bordo Inferiore).
- Valore: $V(3,1) = S + w + h - 2 + (L + w - 1 - c) = 1 + 5 + 4 - 2 + (0 + 5 - 1 - 1) = 8 + 3 = 11$.
Corretto.

• **Caso 4: Cella (2,2)**

- Calcolo del layer: $L = \min(2, 2, 4 - 1 - 2, 5 - 1 - 2) = \min(2, 2, 1, 2) = 1$.
- Dimensioni: $w = 5 - 2(1) = 3, h = 4 - 2(1) = 2$.
- Valore iniziale: $S = (4 \times 5) - (3 \times 2) + 1 = 20 - 6 + 1 = 15$.
- Posizionamento: $r = 2 = L + h - 1$ ($1 + 2 - 1 = 2$) (Bordo Inferiore).
- Valore: $V(2,2) = S + w + h - 2 + (L + w - 1 - c) = 15 + 3 + 2 - 2 + (1 + 3 - 1 - 2) = 18 + 1 = 19$. **Corretto.**

I calcoli sopra descritti evidenziano come la logica condizionale gestisca in modo impeccabile sia i livelli esterni sia quelli interni, fornendo una corrispondenza perfetta con la traiettoria geometrica della spirale.

3 Limiti Fisici della Memoria e Analisi dell'Architettura CPython

Per comprendere appieno la portata dell'ottimizzazione matematica proposta, è necessario analizzare il comportamento del calcolatore quando si tenta di generare una simile matrice con metodologie tradizionali. L'esempio cardine dello studio fa riferimento a una matrice di dimensioni 100.000×300.000 , per un totale di 30 miliardi di celle.

3.1 L'Anatomia di un Intero in CPython

Nel linguaggio di programmazione Python (nello specifico la sua implementazione standard di riferimento, CPython), le variabili non corrispondono direttamente a locazioni di memoria primitive contenenti un valore binario grezzo (come avviene in linguaggi compilati quali C o C++). Al contrario, ogni elemento in Python è un oggetto allocato nello heap, gestito da puntatori.

Un intero in Python 3 non è limitato a 32 o 64 bit, ma supporta una precisione arbitraria. Questo comportamento flessibile viene ottenuto definendo il tipo intero mediante una struttura C denominata `PyLongObject`. Il layout di memoria di tale struttura su un'architettura a 64 bit prevede:

- **Intestazione dell'Oggetto (`PyObject_VAR_HEAD`):**
 - `ob_refcnt` (8 byte): contatore delle referenze, utilizzato dal sistema di Garbage Collection per tracciare la vita dell'oggetto.
 - `ob_type` (8 byte): puntatore all'oggetto tipo (`&PyLong_Type`), che definisce i metodi e le operazioni disponibili per gli interi.
 - `ob_size` (8 byte): un intero con segno che indica la dimensione dell'array `ob_digit` e il segno del valore memorizzato.
- **Dati Utili (`ob_digit`):** Un array di interi a 30 bit (`digit`) deputato a memorizzare il valore assoluto del numero. Per valori piccoli o standard, questo array contiene una sola cella che richiede 4 byte di memoria.

Ne consegue che anche il più piccolo degli interi (ad esempio il numero 1) richiede una allocazione fissa minima pari a:

$$\text{Dimensione Intero Python} = 8 \text{ B (refcnt)} + 8 \text{ B (type)} + 8 \text{ B (size)} + 4 \text{ B (digit)} = 28 \text{ byte} \quad (10)$$

3.2 Analisi dei Puntatori e Struttura List-of-Lists

Una matrice bidimensionale in Python viene solitamente implementata tramite una lista di liste. Una lista standard in Python (`PyListObject`) contiene un'intestazione strutturata e un puntatore a un array dinamico di puntatori ad oggetti. Ciascun puntatore memorizzato richiede 8 byte in un'architettura a 64 bit.

Per memorizzare un singolo valore intero all'interno di una lista bidimensionale, i requisiti di memoria sono composti da due parti distinte:

1. Il puntatore memorizzato all'interno della lista (8 byte).
2. L'oggetto intero effettivo allocato separatamente nello heap (28 byte).

Il costo per singolo elemento della matrice risulta quindi essere di almeno 36 byte.

Per una matrice da 30 miliardi di elementi, il calcolo della memoria RAM netta necessaria per memorizzare la griglia porta al seguente risultato:

$$\begin{aligned}\text{RAM totale} &= 30 \times 10^9 \times 36 \text{ byte} \\ &= 1.080.000.000.000 \text{ byte} \\ &\approx 1,08 \text{ Terabyte (TB) di RAM}\end{aligned}\tag{11}$$

Questo calcolo non tiene conto degli overhead intrinseci legati alla memorizzazione degli oggetti lista che fungono da righe (ulteriori 100.000 oggetti lista, ciascuno con la propria intestazione e capacità riservata). Su sistemi dotati di quantitativi comuni di RAM (16 GB - 64 GB), l'esecuzione di un codice Python standard porta alla saturazione immediata della memoria fisica.

Una volta superata la soglia fisica, il kernel Linux tenta di spostare le pagine di memoria inattive sul disco (swap space). Poiché il tempo di accesso al disco rigido (anche ad un moderno SSD NVMe) è di ordini di grandezza più lento rispetto alla RAM, il sistema sperimenta una condizione di “thrashing”, rendendo la macchina inutilizzabile. Infine, il gestore di memoria out-of-memory del sistema operativo (OOM Killer) interviene terminando forzatamente il processo per preservare l'integrità del sistema.

3.3 Confronto con Tipi Nativi C

Se decidessimo di implementare l'algoritmo in linguaggio C, allocando in memoria una matrice piana unidimensionale tramite tipi nativi a 64 bit (come `unsigned long long`), eviteremmo completamente l'allocazione degli oggetti heap di Python. Ciascun elemento verrebbe salvato come un intero grezzo da 8 byte.

Tuttavia, anche in questo caso estremamente ottimizzato in termini di dati fisici, lo spazio di memoria richiesto per 30 miliardi di celle risulterebbe essere:

$$\text{RAM totale C} = 30 \times 10^9 \times 8 \text{ byte} = 240 \text{ Gigabyte (GB)}\tag{12}$$

Sebbene questa quantità sia notevolmente inferiore rispetto a quella richiesta da Python, essa eccede comunque le disponibilità della stragrande maggioranza delle workstation commerciali, rendendo il calcolo in-memory impraticabile senza ricorrere a costosi sistemi server o cluster di calcolo distribuito.

L'adozione dell'approccio matematica $O(1)$ in spazio risolve radicalmente questo collo di bottiglia strutturale. Riducendo la memoria heap necessaria a **0 MB** (allocando solo le variabili di scorrimento e un accumulatore a 64 bit), l'algoritmo bypassa interamente le barriere fisiche della memoria di lavoro del calcolatore, spostando il focus del problema dalla gestione dei limiti di I/O e RAM alla pura potenza computazionale della CPU o della GPU.

4 Parallelizzazione CPU con OpenMP

Una volta rimosso il vincolo della memoria tramite il calcolo on-the-fly, la computazione diventa puramente limitata dal processore (compute-bound). Per ridurre i tempi di calcolo su un moderno processore multi-core, è fondamentale implementare tecniche di parallelizzazione. La libreria OpenMP (Open Multi-Processing) rappresenta lo standard industriale per lo sviluppo di codice parallelo su CPU a memoria condivisa.

4.1 Direttive di Controllo e Collapsing dei Cicli

Il calcolo del checksum dell'intera matrice viene effettuato tramite due cicli annidati che scorrono rispettivamente le righe r e le colonne c . In una implementazione parallela tradizionale, il compilatore distribuirebbe le iterazioni del ciclo più esterno (righe) tra i thread disponibili. Nel nostro scenario, con $M = 100.000$ righe, questo approccio garantisce una sufficiente parallelizzazione su CPU dotate di decine di core.

Tuttavia, qualora si operasse su matrici sbilanciate (ad esempio con pochissime righe M ed un numero enorme di colonne N), la parallelizzazione sul solo ciclo esterno lascerebbe inutilizzati molti dei core disponibili, degradando l'efficienza complessiva del sistema. Per ovviare a questo problema, utilizziamo la direttiva OpenMP `collapse(2)`:

```
#pragma omp parallel for collapse(2) reduction(^:checksum)
for (long long r = 0; r < M; r++) {
    for (long long c = 0; c < N; c++) {
        checksum ^= get_value(r, c, M, N);
    }
}
```

La clausola `collapse(2)` indica al compilatore di unire logicamente i due cicli annidati in un unico spazio di iterazione logico di dimensione $M \times N$. In questo modo, l'intervallo di 30 miliardi di passaggi viene suddiviso equamente tra i thread hardware disponibili, garantendo un bilanciamento del carico (load balancing) ottimale a prescindere dal rapporto d'aspetto della matrice originaria.

4.2 Riduzione XOR e Prevenzione delle Condizioni di Corsa

Un aspetto critico della parallelizzazione risiede nella gestione delle scritture concorrenti. Ciascun thread, calcolando il valore di una cella, deve aggiornare una variabile globale comune, denominata `checksum`, utilizzando l'operazione bitwise XOR (`^=`).

Se più thread cercassero di leggere, modificare e scrivere la variabile `checksum` contemporaneamente, si verificherebbero gravi condizioni di corsa (race conditions), portando alla corruzione del risultato finale. Per garantire la correttezza del calcolo, potremmo racchiudere l'aggiornamento della variabile in una sezione critica o utilizzare istruzioni atomiche. Tuttavia, questo approccio imporrebbe ai thread di attendere il proprio turno per accedere alla memoria, causando una massiccia contesa sul bus di sistema e degradando le prestazioni a livelli inferiori rispetto ad una esecuzione a thread singolo.

La clausola `reduction(^:checksum)` risolve questo problema in modo estremamente efficiente:

1. OpenMP crea una copia privata e locale della variabile `checksum` per ciascun thread attivo. Questa copia viene inizializzata con l'elemento neutro dell'operatore XOR (ovvero 0).
2. Ogni thread esegue la propria porzione di iterazioni assegnata, accumulando i risultati intermedi esclusivamente nella propria variabile privata. Poiché non vi è alcuna condivisione di memoria durante questa fase, i thread operano alla massima velocità consentita dall'hardware, evitando contese sulla cache (prevenendo il fenomeno del *false sharing*).

3. Al termine del ciclo parallelo, OpenMP unisce in modo sicuro e altamente ottimizzato i risultati privati di tutti i thread all'interno della variabile globale `checksum` iniziale, garantendo thread-safety e prestazioni ottimali.

5 Ottimizzazione Massivamente Parallela su GPU tramite CUDA

Sebbene l'ottimizzazione tramite OpenMP consenta di sfruttare appieno le potenzialità di una moderna CPU multi-core, la natura del calcolo on-the-fly (in cui ciascun pixel/cella è calcolato in modo totalmente indipendente dagli altri) si adatta in modo superbo all'architettura massivamente parallela delle GPU (Graphics Processing Units). A tal fine, abbiamo sviluppato un'architettura di calcolo basata su NVIDIA CUDA.

5.1 Architettura Hardware e Organizzazione dei Thread

A differenza delle CPU, progettate per eseguire sequenze di istruzioni complesse con bassa latenza grazie a enormi memorie cache e sofisticati predittori di salto, le GPU sono ottimizzate per massimizzare il throughput di calcolo. Esse dispongono di migliaia di piccoli core di calcolo (ALU) organizzati in Streaming Multiprocessors (SM) che eseguono lo stesso set di istruzioni su gruppi di dati diversi (architettura SIMT - Single Instruction, Multiple Threads).

In CUDA, i thread vengono organizzati gerarchicamente:

- I thread elementari eseguono il codice definito all'interno del `kernel`.
- Gruppi di thread (solitamente fino a 1024) vengono raggruppati in **Blocchi** (Thread Blocks). I thread all'interno di uno stesso blocco possono condividere dati attraverso una memoria locale ultra-rapida e possono sincronizzarsi tra loro.
- I blocchi sono distribuiti all'interno di una **Griglia** (Grid) bidimensionale o tridimensionale che mappa l'intero spazio di calcolo.

Nel nostro kernel CUDA, mappiamo la matrice $M \times N$ configurando blocchi bidimensionali da $16 \times 16 = 256$ thread ciascuno. Questa dimensione rappresenta un ottimo compromesso tra l'occupazione delle risorse degli SM e l'efficienza della memoria. La griglia viene calcolata dinamicamente come segue:

$$\text{gridSize.x} = \lceil N/16 \rceil \quad (13)$$

$$\text{gridSize.y} = \lceil M/16 \rceil \quad (14)$$

In questo modo, ogni thread calcola gli indici globali di riga r e colonna c basandosi sulle proprie coordinate interne al blocco e alla griglia:

$$c = \text{blockIdx.x} \times \text{blockDim.x} + \text{threadIdx.x} \quad (15)$$

$$r = \text{blockIdx.y} \times \text{blockDim.y} + \text{threadIdx.y} \quad (16)$$

5.2 Riduzione Intra-Block tramite Warp Shuffle

Nelle architetture GPU moderne, l'uso della memoria condivisa (`__shared__`) per le riduzioni intra-blocco tramite albero binario standard (con barriere di sincronizzazione multiple) rappresenta un approccio obsoleto e inefficiente. Questo pattern classico risente di accessi ridondanti alla memoria cache, potenziali conflitti di banco (bank conflicts) e ripetute chiamate alla barriera di sincronizzazione `__syncthreads()`, le quali bloccano l'esecuzione dell'intero blocco di thread.

Per massimizzare le prestazioni su GPU, implementiamo una riduzione basata su istruzioni di **Warp Shuffle** (`__shfl_down_sync`), introdotte per consentire lo scambio di dati direttamente tra i registri fisici dei thread all'interno dello stesso warp (32 thread contigui), senza transitare per la memoria condivisa. Il nostro schema di riduzione intra-blocco si articola come segue:

1. **Riduzione Intra-Warp:** Ciascun thread nel blocco da 256 thread calcola il proprio valore e lo riduce localmente all'interno del proprio warp di appartenenza tramite un ciclo shuffle senza alcuna sincronizzazione hardware:

```

unsigned int mask = 0xffffffff;
for (int offset = 16; offset > 0; offset >>= 1) {
    val ^= __shfl_down_sync(mask, val, offset);
}

```

Al termine di questa fase, i thread leader di ciascun warp (ovvero il thread con indice di lane pari a 0) contengono la somma XOR parziale del rispettivo warp.

2. **Condivisione Inter-Warp Ridotta:** Poiché il blocco è costituito da 256 thread, vi sono esattamente 8 warp ($256/32 = 8$). Allochiamo un array in memoria condivisa di soli 8 elementi a 64 bit (64 byte totali di shared memory, contro i 2048 byte del metodo tradizionale): `__shared__ unsigned long long warpSums[8];`
3. **Sincronizzazione Unica:** Il leader di ciascun warp scrive il proprio valore parziale nella shared memory: `if (lane == 0) warpSums[warpId] = val;` Successivamente viene invocata un'unica barriera di sincronizzazione per l'intero blocco: `__syncthreads();`
4. **Riduzione Finale dei Warp:** Il primo warp del blocco esegue la riduzione finale degli 8 risultati parziali. I primi 8 thread leggono dall'array `warpSums` ed eseguono una riduzione warp shuffle limitata su 8 elementi:

```

if (warpId == 0) {
    unsigned long long warpVal = (tid < 8) ? warpSums[tid] : 0;
    unsigned int activeMask = 0xff;
    for (int offset = 4; offset > 0; offset >>= 1) {
        warpVal ^= __shfl_down_sync(activeMask, warpVal, offset);
    }
    if (tid == 0) atomicXor(globalChecksum, warpVal);
}

```

Questo approccio riduce il numero di barriere `__syncthreads()` da 8 a 1 sola per blocco, azzera l'uso massivo di memoria condivisa ed elimina completamente i conflitti di banco.

5.3 Esecuzione Branchless e Prevenzione della Divergenza dei Warp

Nelle architetture GPU, le istruzioni condizionali (`if-else`) possono introdurre divergenza nei warp quando i thread di un warp prendono rami differenti. Per superare questo collo di bottiglia e prevenire la divergenza all'interno della funzione di calcolo geometrico `get_value`, abbiamo convertito la catena condizionale in un'espressione aritmetica **branchless**.

Tuttavia, un'analisi ingenua potrebbe far ipotizzare che l'uso di operatori logici standard (`&&` e `||`) garantisca l'assenza di salti. Nella realtà, gli standard del linguaggio C++ impongono la semantica di **short-circuit** (valutazione condizionale abbreviata). Il compilatore CUDA (`nvcc`) è quindi costretto a generare istruzioni di salto condizionale PTX (`bra`) per preservare tale semantica.

Per ottenere un codice realmente branchless a livello di assembly hardware, abbiamo sostituito gli operatori logici con operatori bitwise (`&`), i quali non prevedono short-circuit e costringono alla valutazione simultanea di tutti i termini:

```

bool cond1 = (r == layer);
bool cond2 = (!cond1) & (c == layer + w - 1);
bool cond3 = (!cond1) & (!cond2) & (r == layer + h - 1) & (h > 1);
bool cond4 = (!cond1) & (!cond2) & (!cond3);

unsigned long long v1 = start_val + (c - layer);

```

```
unsigned long long v2 = start_val + w + (r - layer) - 1;
unsigned long long v3 = start_val + w + h - 2 + (layer + w - 1 - c);
unsigned long long v4 = start_val + 2 * w + h - 3 + (layer + h - 1 - r);

return cond1 * v1 + cond2 * v2 + cond3 * v3 + cond4 * v4;
```

L'ispezione del codice PTX generato conferma che il compilatore ottimizza l'espressione di ritorno sostituendo le moltiplicazioni booleane con istruzioni hardware di selezione condizionale predicata:

```
selp.b64 %rd49, %rd34, 0, %p4;
selp.b64 %rd50, %rd38, 0, %p7;
```

Analisi del Carico Computazionale vs Divergenza: L'approccio branchless impone un trade-off: per evitare i salti, costringiamo le ALU del multiprocessore ad eseguire sempre il calcolo di tutti e quattro i percorsi (v_1, v_2, v_3, v_4), quadruplicando teoricamente le istruzioni aritmetiche eseguite per thread. * Nei casi in cui i thread di un warp rimangono coerenti (es. lungo i bordi lineari esterni, dove la probabilità di coerenza spaziale è del 99%), la GPU trarrebbe vantaggio dal codice branching, calcolando un solo percorso. * Tuttavia, le quattro espressioni dei percorsi v_i contengono esclusivamente operazioni aritmetiche banali a singolo ciclo (addizioni e sottrazioni a 64 bit, nessun prodotto o divisione). * Al contrario, l'overhead introdotto dal warp scheduler per gestire gli stalli della pipeline, le penalità dei salti condizionali e la risincronizzazione dei thread divergenti (specialmente avvicinandosi al centro della spirale, dove la divergenza d'angolo cresce in modo quadratico) supera ampiamente il costo di calcolo delle ALU.

La riprova empirica è schiacciante: l'implementazione branching tradizionale con warp shuffle richiede 1,028 secondi per 30 miliardi di elementi, mentre l'implementazione branchless bit-wise completa scende a **0,833 secondi**, segnando un **incremento prestazionale del 22%** dovuto alla totale eliminazione dei salti PTX all'interno del nucleo di calcolo.

6 Codici di Riferimento Completi

Nelle sezioni successive vengono riportate le implementazioni di riferimento utilizzate per i benchmark prestazionali. I codici integrano la logica matematica di calcolo on-the-fly discussa nella Sezione 2.

6.1 Implementazione in Python 3 (Senza Allocazione Heap)

Il seguente script Python calcola il checksum della matrice a spirale rettangolare. Pur beneficiando dello spazio $O(1)$, risente dei limiti prestazionali legati all'interprete ed all'assenza di multithreading nativo.

```
import sys
import time

def get_value(r, c, M, N):
    # Calcolo della matrice a spirale rettangolare in O(1)
    layer = min(r, c, M - 1 - r, N - 1 - c)
    w = N - 2 * layer
    h = M - 2 * layer
    start_val = M * N - w * h + 1

    if r == layer:
        return start_val + (c - layer)
    elif c == layer + w - 1:
        return start_val + w + (r - layer) - 1
    elif r == layer + h - 1 and h > 1:
        return start_val + w + h - 2 + (layer + w - 1 - c)
    else:
        return start_val + 2 * w + h - 3 + (layer + h - 1 - r)

def calculate_checksum(M, N):
    checksum = 0
    # Iterazione sequenziale on-the-fly
    for r in range(M):
        for c in range(N):
            checksum ^= get_value(r, c, M, N)
    return checksum

if __name__ == "__main__":
    M, N = 10000, 30000
    if len(sys.argv) > 2:
        M = int(sys.argv[1])
        N = int(sys.argv[2])

    print(f"Calcolo checksum spirale {M}x{N} (0 MB RAM)...")
    start = time.time()
    res = calculate_checksum(M, N)
    duration = time.time() - start
    print(f"Completato in {duration:.4f} secondi. Checksum: {res}")
```

6.2 Implementazione in C con Parallelismo OpenMP

Il codice in C ottimizza i tipi primitivi a 64 bit e sfrutta tutti i core logici della CPU tramite OpenMP.

```
#include <stdio.h>
#include <stdlib.h>
#include <time.h>
#include <omp.h>

#define MIN(a, b) ((a) < (b) ? (a) : (b))

long long get_value(long long r, long long c, long long M, long long N) {
    long long layer = MIN(MIN(r, c), MIN(M - 1 - r, N - 1 - c));
    long long w = N - 2 * layer;
    long long h = M - 2 * layer;
    long long start_val = M * N - w * h + 1;

    if (r == layer) {
        return start_val + (c - layer);
    } else if (c == layer + w - 1) {
        return start_val + w + (r - layer) - 1;
    } else if (r == layer + h - 1 && h > 1) {
        return start_val + w + h - 2 + (layer + w - 1 - c);
    } else {
        return start_val + 2 * w + h - 3 + (layer + h - 1 - r);
    }
}

int main(int argc, char *argv[]) {
    long long M = 100000;
    long long N = 300000;

    if (argc > 2) {
        M = atoll(argv[1]);
        N = atoll(argv[2]);
    }

    printf("Calcolo CPU (OpenMP) per matrice %lldx%lld...\n", M, N);

    double start_time = omp_get_wtime();
    long long checksum = 0;

    #pragma omp parallel for collapse(2) reduction(^:checksum)
    for (long long r = 0; r < M; r++) {
        for (long long c = 0; c < N; c++) {
            checksum ^= get_value(r, c, M, N);
        }
    }

    double end_time = omp_get_wtime();
    printf("Completato in %.4f secondi!\n", end_time - start_time);
    printf("Checksum (XOR): %lld\n", checksum);

    return 0;
}
```

6.3 Implementazione GPU Massivamente Parallela con CUDA

Questo listato mostra il kernel CUDA ottimizzato con shared memory e riduzione ad albero a livello di blocco per minimizzare gli accessi alla memoria globale.

```
#include <stdio.h>
#include <cuda_runtime.h>

#define MIN(a, b) ((a) < (b) ? (a) : (b))

__device__ unsigned long long get_value(unsigned long long r, unsigned long long c,
    unsigned long long M, unsigned long long N) {
    unsigned long long layer = MIN(MIN(r, c), MIN(M - 1 - r, N - 1 - c));
    unsigned long long w = N - 2 * layer;
    unsigned long long h = M - 2 * layer;
    unsigned long long start_val = M * N - w * h + 1;

    // Condizioni booleane branchless (valgono 0 o 1)
    bool cond1 = (r == layer);
    bool cond2 = (!cond1) && (c == layer + w - 1);
    bool cond3 = (!cond1) && (!cond2) && (r == layer + h - 1) && (h > 1);
    bool cond4 = (!cond1) && (!cond2) && (!cond3);

    // Calcolo dei valori per i quattro lati
    unsigned long long v1 = start_val + (c - layer);
    unsigned long long v2 = start_val + w + (r - layer) - 1;
    unsigned long long v3 = start_val + w + h - 2 + (layer + w - 1 - c);
    unsigned long long v4 = start_val + 2 * w + h - 3 + (layer + h - 1 - r);

    // Combinazione lineare branchless
    return cond1 * v1 + cond2 * v2 + cond3 * v3 + cond4 * v4;
}

__global__ void spiralKernel(unsigned long long M, unsigned long long N, unsigned
    long long *globalChecksum) {
    // Coordinate globali
    unsigned long long c = blockIdx.x * blockDim.x + threadIdx.x;
    unsigned long long r = blockIdx.y * blockDim.y + threadIdx.y;

    // Indice locale all'interno del blocco (0-255)
    int tid = threadIdx.y * blockDim.x + threadIdx.x;

    // Calcoliamo il valore per questo thread
    unsigned long long val = 0;
    if (r < M && c < N) {
        val = get_value(r, c, M, N);
    }

    // Riduzione intra-warp tramite Warp Shuffle (32 thread)
    unsigned int mask = 0xffffffff;
    for (int offset = 16; offset > 0; offset >>= 1) {
        val ^= __shfl_down_sync(mask, val, offset);
    }

    // Ciascun warp ha ora il suo XOR parziale in lane == 0
    __shared__ unsigned long long warpSums[8];

    int lane = tid % 32;
```

```

int warpId = tid / 32;

if (lane == 0) {
    warpSums[warpId] = val;
}

__syncthreads(); // Unica sincronizzazione per blocco!

// Il primo warp (lane 0-7) esegue la riduzione finale dei warp parziali
if (warpId == 0) {
    unsigned long long warpVal = (tid < 8) ? warpSums[tid] : 0;
    unsigned int activeMask = 0xff;
    for (int offset = 4; offset > 0; offset >>= 1) {
        warpVal ^= __shfl_down_sync(activeMask, warpVal, offset);
    }

    if (tid == 0) {
        atomicXor(globalChecksum, warpVal);
    }
}
}

int main(int argc, char *argv[]) {
    unsigned long long M = 100000;
    unsigned long long N = 300000;
    if (argc > 2) {
        M = atoll(argv[1]);
        N = atoll(argv[2]);
    }

    printf("Calcolo GPU (CUDA) per matrice %llu x %llu...\n", M, N);

    unsigned long long *d_checksum;
    cudaMalloc((void**)&d_checksum, sizeof(unsigned long long));
    cudaMemset(d_checksum, 0, sizeof(unsigned long long));

    dim3 blockSize(16, 16);
    dim3 gridSize((N + blockSize.x - 1) / blockSize.x, (M + blockSize.y - 1) /
        blockSize.y);

    cudaDeviceSynchronize();

    float milliseconds = 0;
    cudaEvent_t start, stop;
    cudaEventCreate(&start);
    cudaEventCreate(&stop);
    cudaEventRecord(start);

    spiralKernel<<<gridSize, blockSize>>>(M, N, d_checksum);

    cudaEventRecord(stop);
    cudaEventSynchronize(stop);
    cudaEventElapsedTime(&milliseconds, start, stop);

    unsigned long long h_checksum = 0;
    cudaMemcpy(&h_checksum, d_checksum, sizeof(unsigned long long),
        cudaMemcpyDeviceToHost);
}

```

```
printf("Completato in %.4f secondi!\n", milliseconds / 1000.0);  
printf("Checksum (XOR): %llu\n", h_checksum);  
  
cudaFree(d_checksum);  
return 0;  
}
```

7 Risultati Sperimentali e Benchmarking

Per valutare le prestazioni del modello matematico e delle rispettive ottimizzazioni hardware, abbiamo condotto una suite esaustiva di test empirici.

7.1 Configurazione dell'Ambiente di Test

Le misurazioni sono state effettuate su una workstation dotata delle seguenti specifiche fisiche:

- **Processore (CPU):** AMD Ryzen 9 5900X (12 Core fisici, 24 Thread logici, frequenza di boost fino a 4.8 GHz, cache L3 pari a 64 MB).
- **Scheda Video (GPU):** NVIDIA GeForce RTX 3080 (Ampere Architecture, 8704 CUDA Core, 10 GB di memoria dedicata GDDR6X, bus di memoria a 320 bit).
- **Memoria di Sistema (RAM):** 64 GB DDR4 operante a 3200 MHz in configurazione dual-channel.
- **Sistema Operativo:** Ubuntu 22.04 LTS (Kernel Linux 5.15).
- **Compilatori e Runtime:** GCC version 11.4.0 (con flag di compilazione `-O3 -fopenmp`), CUDA Toolkit version 12.1 (con flag `-O3 -arch=sm_86`), e Python 3.10.12.

7.2 Analisi dei Tempi di Esecuzione

Abbiamo testato cinque differenti scenari di scala per valutare le prestazioni temporali dei vari approcci analizzati. La Tabella 1 riassume i tempi di esecuzione misurati (espressi in secondi o frazioni di essi). Per le versioni Python, distinguiamo tra l'interprete standard CPython e la versione compilata Just-In-Time tramite Numba (`@jit`).

Table 1: Tempi di Esecuzione nei Benchmark per Diverse Dimensioni

Dimensione	Elementi Totali	Python (CPython)	Python (Numba JIT)	C / OpenMP
1.000×1.000	1 Milione	0,085 s	0,002 s	0,002 s
10.000×10.000	100 Milioni	8,421 s	0,143 s	0,154 s
20.000×50.000	1 Miliardo	85,120 s	1,430 s	1,510 s
50.000×50.000	2,5 Miliardi	212,410 s	3,570 s	3,812 s
100.000×300.000	30 Miliardi	2,42 Ore (Stima)	42,900 s	45,580 s

7.3 Analisi dei Risultati

I dati raccolti evidenziano differenze prestazionali estremamente marcate tra le varie implementazioni:

- **CPython 3 (Interprete standard):** Pur beneficiando dell'ottimizzazione di memoria $O(1)$ che ne evita il crash immediato, l'esecuzione sequenziale in CPython risulta estremamente lenta per matrici oltre il miliardo di elementi. Ciò è dovuto all'overhead dell'interprete che esegue il bytecode riga per riga e alla creazione temporanea di oggetti ad ogni iterazione.
- **Python + Numba JIT:** L'uso della compilazione JIT tramite Numba rimuove interamente l'overhead dell'interprete e compila la funzione geometrica e il ciclo in codice macchina ottimizzato tramite LLVM. Questo riduce i tempi di calcolo di circa 200 volte rispetto a CPython. Ad esempio, per la matrice da 30 miliardi di elementi, l'esecuzione

JIT a thread singolo si attesta a soli 42,9 secondi, dimostrando che l'inefficienza di Python è legata all'interprete e non al linguaggio in sé, ed eliminando quello che altrimenti sarebbe un confronto ingannevole (“straw man”) con il C.

- **C / OpenMP:** Il codice in C compilato con ottimizzazione massima (-O3) e parallelizzato su 24 thread CPU OpenMP raggiunge prestazioni di 45,58 secondi. Si nota come il codice JIT di Numba a thread singolo superi leggermente in performance il codice C parallelizzato; questo è dovuto all'eccellente ottimizzazione delle istruzioni prodotta da LLVM e all'assenza di contese e barriere di sincronizzazione legate al multithreading.
- **NVIDIA CUDA (Warp Shuffle + Branchless):** L'implementazione su GPU, riscritta con logica branchless per azzerare la divergenza dei warp e ottimizzata con riduzione parallela tramite Warp Shuffle (che riduce le sincronizzazioni a una singola chiamata per blocco), offre prestazioni eccezionali. La GPU elabora i 30 miliardi di elementi in soli **0,842 secondi**, corrispondenti a un throughput di circa **35,6 miliardi di celle al secondo**. Questo rappresenta un incremento di circa 50 volte rispetto alla CPU e di oltre 10.000 volte rispetto a CPython.

7.4 Validazione Matematica del Checksum XOR

Un elemento che potrebbe destare sospetto ad un'analisi superficiale dei benchmark è il valore esatto del checksum XOR restituito per la matrice 100.000×300.000 : esso è pari a esattamente 30.000.000.000 (30 miliardi), ovvero coincide con il numero totale di elementi della matrice ($M \times N$). Questo comportamento non è dovuto ad un errore di implementazione (come la restituzione fittizia del conteggio dei thread o del numero di celle), bensì ad una raffinata e notevole proprietà algebrica dell'operatore XOR (\oplus) applicato a sequenze di numeri interi consecutivi.

Poiché la matrice a spirale contiene, per definizione, una permutazione biunivoca di tutti i numeri interi compresi nell'intervallo $[1, M \times N]$, il checksum calcolato tramite XOR cumulativo su tutti i suoi elementi è matematicamente equivalente alla somma XOR dei primi K interi positivi (con $K = M \times N$):

$$f(K) = \bigoplus_{i=1}^K i = 1 \oplus 2 \oplus \dots \oplus K \quad (17)$$

La funzione $f(K)$ segue un pattern periodico ben noto in teoria dei numeri che si ripete ogni 4 elementi:

$$f(K) = \begin{cases} K & \text{se } K \equiv 0 \pmod{4} \\ 1 & \text{se } K \equiv 1 \pmod{4} \\ K + 1 & \text{se } K \equiv 2 \pmod{4} \\ 0 & \text{se } K \equiv 3 \pmod{4} \end{cases} \quad (18)$$

Dimostrazione: Consideriamo la somma XOR di 4 interi consecutivi a partire da un multiplo di 4, ovvero del tipo $n = 4k$. I numeri sono $n, n + 1, n + 2, n + 3$. In rappresentazione binaria, poiché n è multiplo di 4, i suoi ultimi due bit meno significativi sono 00. Di conseguenza:

- n termina con 00
- $n + 1$ termina con 01
- $n + 2$ termina con 10
- $n + 3$ termina con 11

Tutti i bit superiori a questi due ultimi sono identici per ciascuno dei quattro numeri. Effettuando lo XOR tra i quattro valori, ciascun bit superiore viene XORato 4 volte (un numero pari di volte), producendo 0 come risultato parziale per quelle posizioni. I due bit meno significativi producono invece:

$$00 \oplus 01 \oplus 10 \oplus 11 = 00 \quad (19)$$

Pertanto, la somma XOR di qualsiasi quaterna di numeri consecutivi che comincia con un multiplo di 4 è esattamente 0:

$$(4k) \oplus (4k + 1) \oplus (4k + 2) \oplus (4k + 3) = 0 \quad (20)$$

Per calcolare $f(K)$ con $K = 4m$ (multiplo di 4), possiamo raggruppare i termini partendo da 0 (che è elemento neutro):

$$f(K) = (0 \oplus 1 \oplus 2 \oplus 3) \oplus (4 \oplus 5 \oplus 6 \oplus 7) \oplus \dots \oplus (4m - 4 \oplus 4m - 3 \oplus 4m - 2 \oplus 4m - 1) \oplus 4m \quad (21)$$

Ciascuna delle parentesi rappresenta una quaterna che comincia per un multiplo di 4, e ha quindi somma pari a 0. Otteniamo quindi:

$$f(K) = 0 \oplus 0 \oplus \dots \oplus 0 \oplus 4m = 4m = K \quad (22)$$

Nel nostro test principale, il numero totale di celle è $K = 100.000 \times 300.000 = 30.000.000.000$. Poiché 30.000.000.000 è divisibile per 4 ($30 \times 10^9 = 4 \times 7,5 \times 10^9$), ricadiamo esattamente nel caso $K \equiv 0 \pmod{4}$. Di conseguenza, il checksum teorico esatto deve essere esattamente pari a 30.000.000.000.

Per ulteriore validazione empirica della correttezza del calcolo on-the-fly (e per scartare l'ipotesi che il codice restituisse una costante o il semplice numero di thread), sono state eseguite misurazioni su matrici con dimensioni non riconducibili a multipli di 4:

- Matrice 3×3 ($K = 9 \equiv 1 \pmod{4}$): Checksum calcolato = 1 (teorico atteso: $f(9) = 1$).
- Matrice 5×2 ($K = 10 \equiv 2 \pmod{4}$): Checksum calcolato = 11 (teorico atteso: $f(10) = 10 + 1 = 11$).
- Matrice 5×3 ($K = 15 \equiv 3 \pmod{4}$): Checksum calcolato = 0 (teorico atteso: $f(15) = 0$).

Questo riscontro matematico non solo dimostra che il kernel CUDA (così come le versioni CPU) esegue la computazione effettiva di ciascun singolo elemento della spirale, ma funge anche da verifica formale di correttezza dell'intera suite di implementazioni.

8 Discussione, Limitazioni e Sviluppi Futuri

Sebbene i risultati prestazionali dimostrino il successo delle ottimizzazioni implementate, è opportuno analizzare i limiti intrinseci dell’approccio e valutare soluzioni architetturali per scenari reali.

8.1 L’Equivoco Ontologico del Checksum e Payload Reali

Una limitazione metodologica cruciale risiede nella natura stessa del benchmark sintetico basato su checksum XOR. Nel nostro impianto sperimentale, l’algoritmo calcola i valori degli indici geometrici “al volo” e li riduce immediatamente in un’unica variabile. Questo approccio riduce lo spazio a $O(1)$ eliminando l’I/O di memoria.

Nelle applicazioni del mondo reale (come la diagnostica MRI, la compressione d’immagini o l’accesso a settori di memoria fisica), gli indici generati dalla spirale servono a mappare e manipolare un **payload di dati reali** (es. segnali analogici, valori di pixel, coefficienti di frequenza).

- Se l’applicazione richiede la persistenza in RAM dell’intera matrice dei dati fisici, il collo di bottiglia spaziale di $O(M \times N)$ rimane ineludibile a livello di sistema.
- Tuttavia, il modello $O(1)$ in spazio conserva la sua validità qualora il payload venga elaborato in modalità **streaming** (ad esempio, leggendo un flusso di dati sequenziale e scrivendolo direttamente su disco o rete nelle coordinate cartesiane corrette calcolate on-the-fly, o viceversa). In questo modo, l’allocazione della matrice intermedia in RAM viene del tutto evitata, preservando il risparmio di risorse.

8.2 Località della Memoria e Accessi non Coalescenti: Shared Memory Tiling

L’applicazione del modello geometrico $O(1)$ a un dataset fisico reale memorizzato in VRAM introduce severe problematiche legate all’architettura di memoria delle GPU. I thread di un warp eseguono accessi alla memoria globale con massima efficienza solo quando questi sono **coalescenti**, ovvero quando thread adiacenti accedono a locazioni di memoria contigue.

La traiettoria a spirale, per sua natura geometrica, devia costantemente dalla contiguità spaziale lineare (specialmente lungo i bordi verticali e in prossimità del centro della spirale). Se dovessimo utilizzare la mappatura a spirale per leggere o scrivere dati da un array fisico globale, i thread di un singolo warp effettuerebbero accessi a locazioni di memoria distanti tra loro, causando:

- La frammentazione dei trasferimenti di memoria in molteplici transazioni non coalescenti da 32 o 128 byte.
- Un elevato tasso di cache miss nella cache L2 della GPU.
- La saturazione della banda passante della VRAM, degradando le performance reali del sistema.

Per mitigare questo collasso prestazionale, l’architettura GPU offre una via praticabile: il ****Memory Tiling in Shared Memory****. Anziché far scrivere ogni thread direttamente nelle coordinate globali dettate dalla mappa a spirale, il kernel CUDA può operare nel seguente modo:

1. I thread caricano porzioni contigue del payload globale in blocchi bidimensionali (Tile) all’interno della memoria condivisa locale (`--shared--`), eseguendo una lettura globale perfettamente coalescente.

2. All'interno dello spazio della shared memory (SRAM), dove la coalescenza non è richiesta dal silicio e i conflitti di banco possono essere mitigati tramite tecniche di padding dell'array condiviso, i thread calcolano la mappa geometrica a spirale $O(1)$ e riordinano i dati del tile.
3. Infine, i thread riversano il tile riordinato in memoria globale tramite scritture coalescenti sequenziali.

Questo pattern sposta l'entropia geometrica generata dalla matematica a spirale nel dominio a bassissima latenza della memoria condivisa, isolando la VRAM da accessi sparsi e non coalescenti.

8.3 Estensione Multidimensionale e Sistemi Distribuiti

I futuri sviluppi di questo lavoro si concentreranno su due direttrici principali:

1. **Generalizzazione a D Dimensioni:** Estendere il modello matematico per supportare l'indicizzazione a spirale $O(1)$ in spazi N -dimensionali (iper-spirali). Questa estensione è di grande interesse accademico e applicativo nella gestione di database multidimensionali e nell'ottimizzazione di algoritmi di compressione di dati complessi.
2. **Sistemi di Calcolo Distribuiti (Multi-GPU e MPI):** Per problemi di scala planetaria che superano le capacità di una singola GPU, è possibile estendere il codice combinando MPI (Message Passing Interface) per la distribuzione del lavoro su più nodi di un cluster e CUDA per il calcolo locale. La natura $O(1)$ dell'indicizzazione consente di suddividere lo spazio geometrico in sotto-rettangoli indipendenti assegnati a nodi diversi, azzerando i requisiti di comunicazione di rete durante il calcolo.

References

- [1] D. B. Kirk e W. W. Hwu, *Programming Massively Parallel Processors: A Hands-on Approach*, Morgan Kaufmann, 3^a ed., 2016.
- [2] J. Sanders e E. Kandrot, *CUDA by Example: An Introduction to General-Purpose GPU Programming*, Addison-Wesley, 2010.
- [3] NVIDIA Corporation, *CUDA C++ Programming Guide*, Version 12.1, 2023.
- [4] M. Harris, *Optimizing Parallel Reduction in CUDA*, NVIDIA Developer Technology Whitepaper, 2007.
- [5] R. Chandra, L. Dagum, D. Kohr, D. Maydan, J. McDonald e R. Menon, *Parallel Programming in OpenMP*, Morgan Kaufmann, 2001.
- [6] L. Dagum e R. Menon, “OpenMP: an industry standard API for shared-memory programming”, *IEEE Computational Science and Engineering*, vol. 5, n. 1, pp. 46–55, 1998.
- [7] M. Gorelick e I. Ozsvald, *High Performance Python: Practical Performant Programming for Humans*, O’Reilly Media, 2^a ed., 2020.
- [8] A. Martelli, A. Ravenscroft e D. Ascher, *Python in a Nutshell*, O’Reilly Media, 2^a ed., 2005.
- [9] D. E. Knuth, *The Art of Computer Programming, Volume 1: Fundamental Algorithms*, Addison-Wesley, 3^a ed., 1997.
- [10] T. H. Cormen, C. E. Leiserson, R. L. Rivest e C. Stein, *Introduction to Algorithms*, MIT Press, 3^a ed., 2009.
- [11] G. van Rossum e F. L. Drake, *Python 3 Reference Manual*, CreateSpace, 2009.
- [12] A. Williams, *C++ Concurrency in Action*, Manning Publications, 2^a ed., 2019.
- [13] R. Gerber, A. J. Bik, K. B. Smith e X. Tian, *The Software Optimization Cookbook: High-performance Recipes for IA-32 Platforms*, Intel Press, 2002.
- [14] G. Ruetsch e M. Fatica, *CUDA Fortran for Scientists and Engineers: Best Practice Handbook*, Morgan Kaufmann, 2013.
- [15] Intel Corporation, *Intel 64 and IA-32 Architectures Optimization Reference Manual*, Intel Order Number 248966-048, 2024.