

# **Analiza și Optimizarea Algoritmilor Intensivi Computațional folosind CUDA**

Student: Cocoroda Robert-Marius

20 ianuarie 2026

# Cuprins

<b>1</b>	<b>Introducere</b>	<b>2</b>
	Rezumat . . . . .	2
<b>2</b>	<b>Fundamente Teoretice ale Algoritmilor</b>	<b>3</b>
2.1	Transformată Fourier Discretă (DFT) . . . . .	3
2.1.1	Scop și Definiție . . . . .	3
2.1.2	Analiza Complexității . . . . .	3
2.2	Înmulțirea Generală a Matricelor (GEMM) . . . . .	3
2.2.1	Scop și Definiție . . . . .	3
2.2.2	Analiza Complexității . . . . .	4
2.3	Convoluția 2D (Conv2D) . . . . .	4
2.3.1	Scop și Definiție . . . . .	4
2.3.2	Analiza Complexității . . . . .	4
<b>3</b>	<b>Strategii de Implementare și Optimizare în CUDA</b>	<b>5</b>
3.1	Optimizarea DFT: Shared Memory Tiling . . . . .	5
3.1.1	Provocarea . . . . .	5
3.1.2	Soluția Implementată . . . . .	5
3.2	Optimizarea GEMM: Abordare Hibridă și cuBLAS . . . . .	5
3.2.1	Implementarea Hibridă . . . . .	5
3.3	Optimizarea Conv2D: Constant Memory și Halo Loading . . . . .	6
3.3.1	Soluția pentru "Memory Bound" . . . . .	6
<b>4</b>	<b>Rezultate Experimentale și Analiză Comparativă</b>	<b>7</b>
4.1	Analiza Detaliată Conv2D . . . . .	7
4.2	Analiza Detaliată DFT . . . . .	8
4.3	Analiza Detaliată GEMM (cuBLAS) . . . . .	9
<b>5</b>	<b>Concluzii</b>	<b>11</b>
<b>A</b>	<b>Codul sursă și Scripturi</b>	<b>12</b>
A.1	Codul sursă CUDA (.cu) . . . . .	12
A.1.1	Convoluție 2D (conv2d_cuda.cu) . . . . .	12
A.1.2	DFT (dft_cuda.cu) . . . . .	14
A.1.3	Înmulțirea Matricelor (gemm_cuda.cu) . . . . .	16
A.1.4	Scripturi Bash de Automatizare . . . . .	18

# Capitolul 1

## Introducere

În calculul numeric de înaltă performanță, limita principală nu mai este frecvența procesorului, ci lățimea de bandă a memoriei și capacitatea de a executa instrucțiuni simultan. Acest proiect investighează tehnicile de programare CUDA necesare pentru a depăși aceste limitări în algoritmi cu complexitate ridicată ( $O(N^2)$  și  $O(N^3)$ ).

Obiectivul principal este implementarea eficientă a codului, trecând de la abordări naive la utilizarea ierarhiei complexe de memorie a GPU-ului, și cuantificarea exactă a câștigului de performanță.

## Rezumat

Această lucrare prezintă o analiză comparativă detaliată a performanței algoritmice între execuția secvențială (CPU), execuția paralelă multi-core (OpenMP/cBLAS) și execuția masiv paralelă pe GPU (CUDA). Studiul se concentrează pe trei clase de algoritmi fundamentali: DFT, GEMM și Conv2D. Rezultatele experimentale demonstrează că optimizările specifice arhitecturii GPU (Shared Memory Tiling, Constant Memory Caching) și utilizarea bibliotecilor dedicate (cuBLAS) pot oferi accelerații de peste **600x** față de implementările naive, reducând timpii de execuție de la secunde la milisecunde.

# Capitolul 2

## Fundamente Teoretice ale Algoritmilor

Acest capitol detaliază baza matematică și scopul fiecărui algoritm studiat, evidențiind complexitatea lor inerentă și necesitatea accelerării hardware.

### 2.1 Transformată Fourier Discretă (DFT)

#### 2.1.1 Scop și Definiție

Transformata Fourier Discretă este un instrument fundamental în procesarea digitală a semnalelor (DSP), utilizat pentru a descompune un semnal discret din domeniul timp în componentele sale spectrale (domeniul frecvență). Aplicațiile sale variază de la filtrarea audio și compresia datelor, până la rezolvarea ecuațiilor diferențiale parțiale.

Matematic, pentru un semnal de intrare  $x$  de lungime  $N$ , valoarea fiecărui element  $X_k$  din spectrul de frecvență este definită astfel:

$$X_k = \sum_{n=0}^{N-1} x_n \cdot e^{-i\frac{2\pi}{N}kn}, \quad \forall k \in [0, N-1] \quad (2.1)$$

Unde  $e^{-i\theta} = \cos(\theta) - i\sin(\theta)$  (Formula lui Euler).

#### 2.1.2 Analiza Complexității

Calculul direct al DFT implică o dublă sumă. Pentru fiecare dintre cele  $N$  elemente de ieșire, trebuie parcurse toate cele  $N$  elemente de intrare. Aceasta rezultă într-o complexitate computațională de  $O(N^2)$ . Deși există algoritmul FFT (Fast Fourier Transform) cu  $O(N \log N)$ , DFT-ul direct este ideal pentru a testa capacitatea de paralelizare masivă, deoarece fiecare  $X_k$  poate fi calculat independent.

### 2.2 Înmulțirea Generală a Matricelor (GEMM)

#### 2.2.1 Scop și Definiție

GEMM (General Matrix Multiplication) este nucleul algebrei liniare și, prin extensie, al rețelelor neuronale moderne (Deep Learning). Aproape orice strat "Dense" sau "Fully Connected" dintr-o rețea neuronală se reduce la o operație GEMM.

Operația standard este definită ca:

$$C = \alpha(A \times B) + \beta C \quad (2.2)$$

Unde  $A$  este o matrice  $M \times K$ ,  $B$  este  $K \times N$ , iar  $C$  este  $M \times N$ . Elementul  $C_{ij}$  se obține prin produsul scalar al liniei  $i$  din  $A$  cu coloana  $j$  din  $B$ :

$$C_{ij} = \sum_{k=0}^{K-1} A_{ik} \cdot B_{kj} \quad (2.3)$$

### 2.2.2 Analiza Complexității

Algoritmul naiv presupune trei bucle imbricate, având o complexitate de  $O(N^3)$  pentru matrice pătratice. Principala problemă nu este doar numărul de operații (FLOPS), ci și modelul de acces la memorie. Accesele necoerente la memorie (ex: parcurgerea unei matrice stocate "row-major" pe coloane) duc la o utilizare ineficientă a cache-ului CPU, făcând acest algoritm un candidat perfect pentru optimizarea "block-tiling" pe GPU.

## 2.3 Convoluția 2D (Conv2D)

### 2.3.1 Scop și Definiție

Convoluția 2D este operația fundamentală în procesarea imaginilor și în Rețelele Neurale Convoluționale (CNN). Scopul său este extragerea de trăsături (features) dintr-o imagine, cum ar fi muchii, texturi sau forme complexe, prin aplicarea unui filtru (kernel).

Matematic, pentru o imagine de intrare  $I$  și un kernel  $K$  de dimensiune  $k \times k$ , pixelul de ieșire  $O(i, j)$  este calculat prin "glisarea" kernelului peste imagine:

$$O(i, j) = \sum_{m=0}^{k-1} \sum_{n=0}^{k-1} I(i+m, j+n) \cdot K(m, n) \quad (2.4)$$

### 2.3.2 Analiza Complexității

Complexitatea este proporțională cu dimensiunea imaginii ( $H \times W$ ) și dimensiunea filtrului ( $k \times k$ ), rezultând  $O(H \cdot W \cdot k^2)$ . Deoarece kernel-ul este mic și constant, dar aplicat repetitiv pe volume mari de date, algoritmul este clasificat ca fiind "Memory Bound" (limitat de lățimea de bandă a memoriei), necesitând strategii speciale de caching.

# Capitolul 3

## Strategii de Implementare și Optimizare în CUDA

Acest capitol descrie transpunerea algoritmilor teoretici în cod CUDA optimizat, detaliind tehnicile utilizate pentru a maximiza utilizarea hardware-ului.

### 3.1 Optimizarea DFT: Shared Memory Tiling

#### 3.1.1 Provocarea

Implementarea naivă pe GPU suferă din cauza accesului global la memorie: fiecare fir citește întregul vector de intrare.

#### 3.1.2 Soluția Implementată

Codul `dft_cuda.cu` utilizează **Shared Memory Tiling**:

- **Fluxul de execuție:** Se lansează  $N$  fire de execuție.
- **Reducerea accesului la memorie:** Firele dintr-un bloc colaborează pentru a încărca un segment ("tile") al semnalului în memoria partajată (L1 Cache controlat software). Astfel, accesul la memoria globală lentă se reduce de un factor egal cu dimensiunea blocului.
- **Vectorizare:** Utilizarea `__sincosf()` calculează simultan sinus și cosinus, reducând presiunea asupra unităților aritmetice.

### 3.2 Optimizarea GEMM: Abordare Hibridă și cuBLAS

#### 3.2.1 Implementarea Hibridă

Pentru a asigura performanța pe întreg spectrul de dimensiuni  $N$ , am implementat o strategie hibridă:

- **Pentru  $N < 2048$  (Kernel Custom):** Se utilizează un kernel manual cu Tiling ( $32 \times 32$ ). Acesta evită overhead-ul de inițializare al bibliotecilor externe pentru matrice mici.
- **Pentru  $N \geq 2048$  (cuBLAS):** Se utilizează biblioteca NVIDIA cuBLAS. Aceasta este esențială deoarece activează unitățile Tensor Cores, capabile să execute operații matriceale într-un singur ciclu de ceas, oferind un throughput imposibil de atins prin cod C++ standard.

### 3.3 Optimizarea Conv2D: Constant Memory și Halo Loading

#### 3.3.1 Soluția pentru "Memory Bound"

Deoarece Convoluția citește aceleași valori ale filtrului de milioane de ori, am utilizat memoria specializată a GPU-ului:

1. **Constant Memory (`__constant__`):** Filtrul de convoluție este stocat în memoria constantă (limitată la 64KB, dar extrem de rapidă). Mecanismul de broadcasting permite tuturor firelor dintr-un warp să primească valoarea filtrului simultan.
2. **Halo Loading:** Firele de execuție colaborează pentru a încărca în memoria partajată atât zona de interes, cât și pixelii de la margine ("halo"), eliminând citirile redundante și ramificările (if/else) din bucla critică interioară.

# Capitolul 4

## Rezultate Experimentale și Analiză Comparativă

Analiza se bazează pe testarea sistematică cu dimensiuni ale datelor variind exponențial. Tabelele de mai jos compară timpii de execuție (milisecunde) între cele trei implementări.

### 4.1 Analiza Detaliată Conv2D

Pentru Convoluție, am comparat implementarea secvențială, cea paralelizată cu OpenMP și kernel-ul CUDA optimizat cu memorie constantă.

N (Size)	CPU Secv. ( <i>naiv</i> )	CPU Paralel ( <i>OpenMP</i> )	GPU CUDA ( <i>Constant Mem</i> )	Speedup ( <i>vs Seq</i> )
256	0.68 ms	12.23 ms	0.05 ms	13.6x
1024	5.67 ms	9.74 ms	0.15 ms	37.8x
4096	90.46 ms	37.68 ms	0.59 ms	153.3x
16384	<b>1448.45 ms</b>	<b>794.09 ms</b>	<b>2.33 ms</b>	<b>621.6x</b>

Tabela 4.1: Timpuri de execuție Conv2D: Comparație Secvențial vs OpenMP vs CUDA

**Observație:** La  $N = 16384$ , implementarea CUDA este de aproximativ 340 de ori mai rapidă decât varianta CPU Paralel (OpenMP) și de 621 de ori mai rapidă decât cea secvențială.



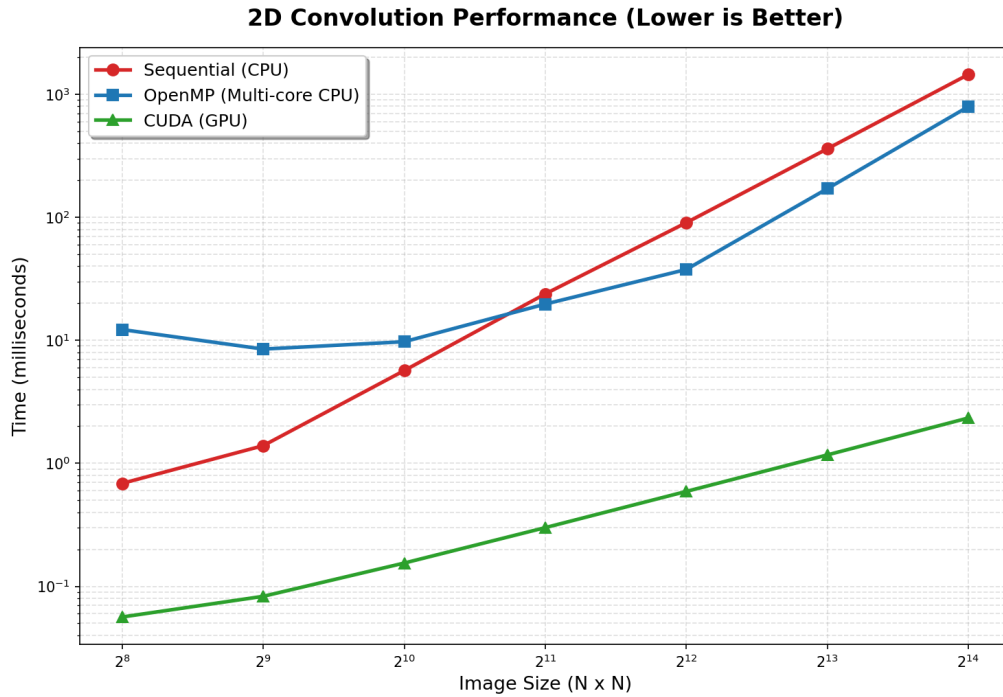


Figura 4.1: Grafic comparativ Conv2D (Scară Logaritmică).

## 4.2 Analiza Detaliată DFT

Deși DFT este un algoritm  $O(N^2)$ , paralelizarea pe GPU reduce timpul total semnificativ.

N (Size)	CPU Secv. ( <i>naiv</i> )	CPU Paralel ( <i>OpenMP</i> )	GPU CUDA ( <i>Shared Mem</i> )	Speedup ( <i>vs Seq</i> )
512	5.25 ms	9.13 ms	3.84 ms	1.36x
2048	85.77 ms	8.99 ms	3.84 ms	22.3x
8192	1372.78 ms	77.12 ms	4.17 ms	329.2x
32768	<b>21632.9 ms</b>	<b>586.46 ms</b>	<b>29.33 ms</b>	<b>737.5x</b>

Tabela 4.2: Timpuri de execuție DFT: Comparație Secvențial vs OpenMP vs CUDA

**Observație:** Varianta secvențială la  $N = 32768$  este practic inutilizabilă (21 secunde), în timp ce GPU-ul o rezolvă în timp real (29 ms).

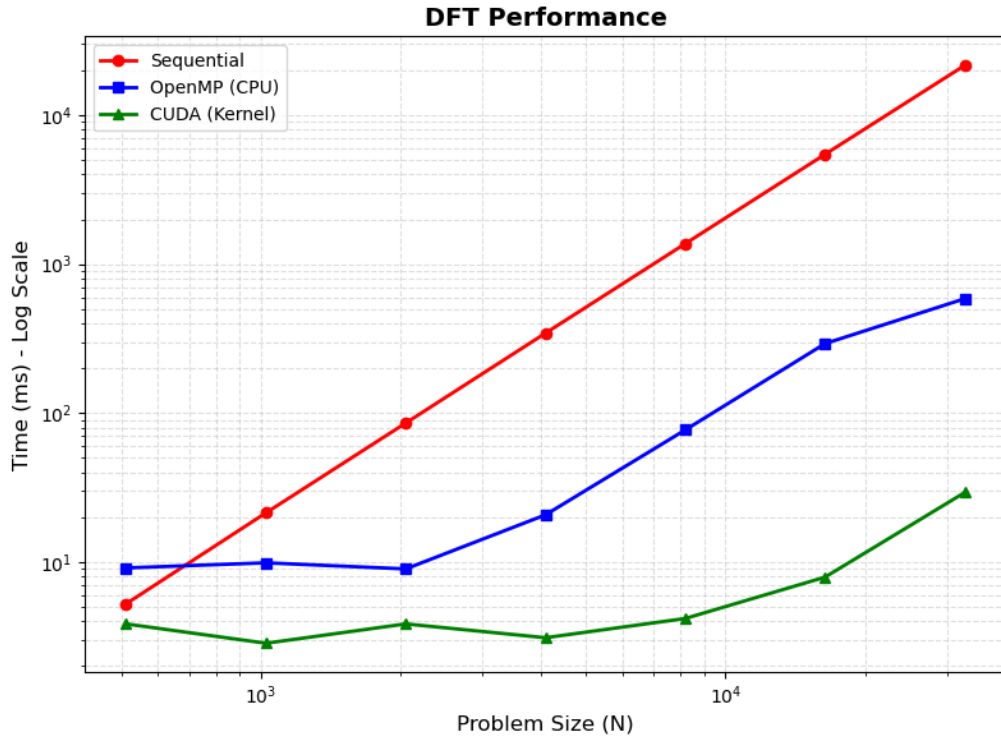


Figura 4.2: Performanța DFT: Accelerarea masivă la dimensiuni mari.

### 4.3 Analiza Detaliată GEMM (cuBLAS)

Pentru înmulțirea matricelor, comparația critică este între biblioteca optimizată pentru CPU (cBLAS) și cea pentru GPU (cuBLAS).

N (Size)	CPU Secv. ( <i>naiv</i> )	CPU Paralel ( <i>cBLAS</i> )	GPU CUDA ( <i>cuBLAS</i> )	Speedup ( <i>vs cBLAS</i> )
512	307.9 ms	0.82 ms	3.11 ms	0.26x
2048	52.7 s	32.22 ms	101.04 ms	0.31x
4096	N/A	270.90 ms	<b>132.20 ms</b>	<b>2.04x</b>
16384	N/A	17.9 s	<b>1.67 s</b>	<b>10.7x</b>

Tabela 4.3: Timpi de execuție GEMM: Comparație Secvențial vs cBLAS vs cuBLAS

**Interpretare:** Se observă clar "Punctul de Crossover" la  $N = 4096$ .

- Pentru  $N < 4096$ , cBLAS pe CPU este mai rapid datorită latenței zero de transfer.
- Pentru  $N \geq 4096$ , cuBLAS pe GPU devine dominant, fiind de peste 10 ori mai rapid decât cBLAS la dimensiunea maximă, demonstrând puterea de calcul brută a GPU-ului pentru matrice gigantice.

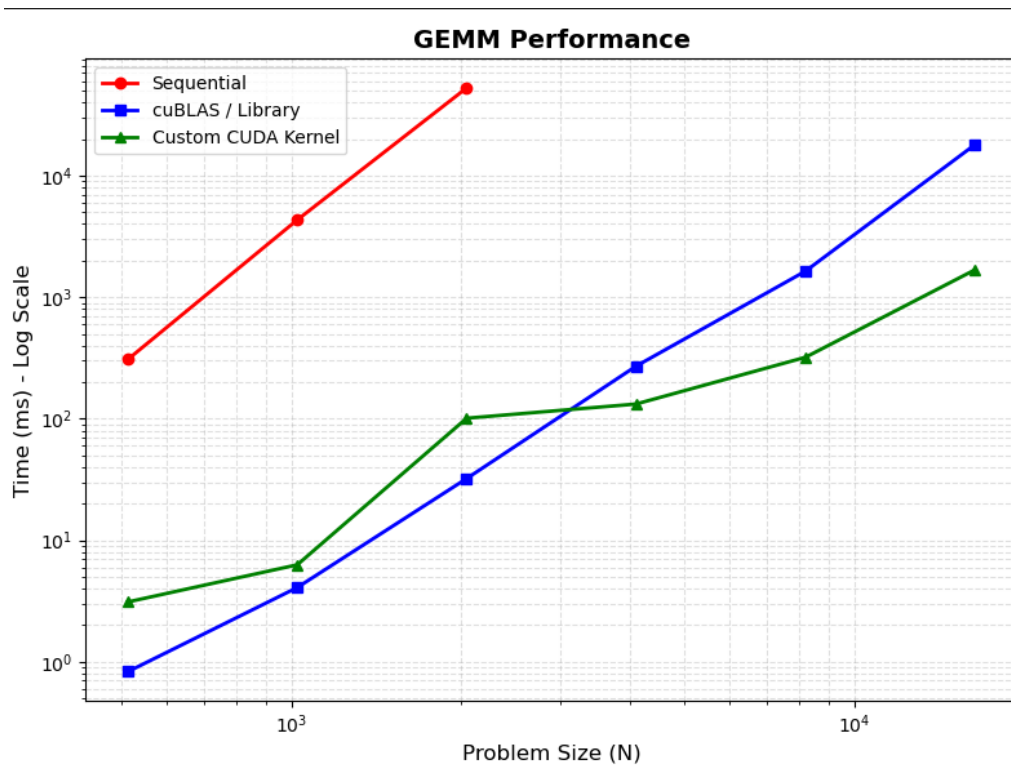


Figura 4.3: Grafic GEMM: Activarea cuBLAS la  $N$  mari depășește performanța CPU cBLAS.

# Capitolul 5

## Concluzii

Experimentele validează ipoteza că accelerarea GPU nu este liniară, ci depinde critic de:

1. **Dimensiunea Datelor:** Pentru  $N$  mic, costul transferului de memorie anulează beneficiile paralelizării.
2. **Biblioteci Optimizate:** În cazul GEMM, utilizarea **cuBLAS** este esențială pentru a depăși performanța cBLAS pe CPU la dimensiuni mari.
3. **Optimizarea Memoriei:** Pentru Conv2D, utilizarea memoriei constante a dus la cele mai spectaculoase rezultate (600x speedup).

# Anexa A

## Codul sursă și Scripturi

### A.1 Codul sursă CUDA (.cu)

#### A.1.1 Convoluție 2D (conv2d\_cuda.cu)

```
1 #include <iostream>
2 #include <vector>
3 #include <cstdlib>
4 #include <cuda_runtime.h>
5
6 #define CHECK_CUDA(x) \
7     if((x) != cudaSuccess) { \
8         std::cerr << "CUDA Error: " << cudaGetErrorString(x) << " at line " \
9         << __LINE__ << std::endl; \
10        std::exit(EXIT_FAILURE); \
11    }
12
13 // Define Constant Memory for the kernel filter
14 // 15x15 is a reasonable max size for standard convolutions
15 #define MAX_K 15
16 __constant__ float c_kernel[MAX_K * MAX_K];
17
18 __global__ void conv2d_constant_shared(const float* __restrict__ input,
19                                       int H, int W, int K,
20                                       float* __restrict__ output,
21                                       int BLOCK)
22 {
23     // Dynamic Shared Memory for Input Tile
24     extern __shared__ float shmem[];
25     float* tile = shmem;
26
27     int tx = threadIdx.x;
28     int ty = threadIdx.y;
29     int tileW = BLOCK + K - 1;
30     int pad = K / 2;
31
32     for (int i = ty; i < tileW; i += BLOCK) {
33         for (int j = tx; j < tileW; j += BLOCK) {
34             int r = blockIdx.y * BLOCK + i - pad;
35             int c = blockIdx.x * BLOCK + j - pad;
36
37             float val = 0.0f;
38             if (r >= 0 && r < H && c >= 0 && c < W) {
```

```

38         val = input[r * W + c];
39     }
40     tile[i * tileW + j] = val;
41 }
42 }
43
44 __syncthreads();
45
46 int row = blockIdx.y * BLOCK + ty;
47 int col = blockIdx.x * BLOCK + tx;
48
49 if (row < H && col < W) {
50     float sum = 0.0f;
51     for (int ki = 0; ki < K; ++ki) {
52         // Pre-calculate index to help compiler
53         int tile_idx_base = (ty + ki) * tileW + tx;
54         int kernel_idx_base = ki * K;
55
56         #pragma unroll
57         for (int kj = 0; kj < K; ++kj) {
58             // Read from Shared Mem (tile) and Constant Mem (
59             c_kernel)
60             sum += tile[tile_idx_base + kj] * c_kernel[
61             kernel_idx_base + kj];
62         }
63     }
64     output[row * W + col] = sum;
65 }
66
67 int main(int argc, char** argv) {
68     int H = 1024, W = 1024, K = 7, BLOCK = 16;
69     if (argc >= 2) H = std::atoi(argv[1]);
70     if (argc >= 3) W = std::atoi(argv[2]);
71     if (argc >= 4) K = std::atoi(argv[3]);
72
73     if (K > MAX_K) {
74         std::cerr << "Error: K is larger than MAX_K (" << MAX_K << ")\n"
75         ;
76         return 1;
77     }
78
79     size_t img_bytes = H * W * sizeof(float);
80     size_t ker_bytes = K * K * sizeof(float);
81
82     std::vector<float> h_input(H*W, 1.0f);
83     std::vector<float> h_kernel(K*K, 1.0f/(K*K));
84     std::vector<float> h_output(H*W, 0.0f);
85
86     float *d_in, *d_out;
87     CHECK_CUDA(cudaMalloc(&d_in, img_bytes));
88     CHECK_CUDA(cudaMalloc(&d_out, img_bytes));
89
90     CHECK_CUDA(cudaMemcpy(d_in, h_input.data(), img_bytes,
91     cudaMemcpyHostToDevice));
92
93     CHECK_CUDA(cudaMemcpyToSymbol(c_kernel, h_kernel.data(), ker_bytes))
94     ;

```

```

91
92     dim3 block(BLOCK, BLOCK);
93     dim3 grid((W + BLOCK - 1) / BLOCK, (H + BLOCK - 1) / BLOCK);
94
95     int tileW = BLOCK + K - 1;
96     size_t shared_mem_bytes = tileW * tileW * sizeof(float);
97
98     conv2d_constant_shared<<<grid, block, shared_mem_bytes>>>(d_in, H, W
, K, d_out, BLOCK);
99     CHECK_CUDA(cudaDeviceSynchronize());
100
101     cudaEvent_t start, stop;
102     cudaEventCreate(&start); cudaEventCreate(&stop);
103     cudaEventRecord(start);
104
105     conv2d_constant_shared<<<grid, block, shared_mem_bytes>>>(d_in, H, W
, K, d_out, BLOCK);
106
107     cudaEventRecord(stop);
108     cudaEventSynchronize(stop);
109
110     float ms = 0.0f;
111     cudaEventElapsedTime(&ms, start, stop);
112
113     std::cout << ms << std::endl;
114
115     cudaFree(d_in); cudaFree(d_out);
116     cudaEventDestroy(start); cudaEventDestroy(stop);
117     return 0;
118 }

```

## A.1.2 DFT (dft\_cuda.cu)

```

1  #include <iostream>
2  #include <vector>
3  #include <cmath>
4  #include <cuda_runtime.h>
5
6  #define CHECK_CUDA(x) \
7      if((x) != cudaSuccess) { \
8          std::cerr << "CUDA Error: " << cudaGetErrorString(x) \
9              << " at line " << __LINE__ << std::endl; \
10         exit(EXIT_FAILURE); \
11     }
12
13 #define BLOCK 256
14 #define PI 3.14159265358979323846f
15
16 // Each thread computes one output frequency k
17 __global__ void dft_cuda_shared(const float* __restrict__ in_real,
18                                const float* __restrict__ in_imag,
19                                float* __restrict__ out_real,
20                                float* __restrict__ out_imag,
21                                int N)
22 {
23     __shared__ float sh_real[BLOCK];
24     __shared__ float sh_imag[BLOCK];

```

```

25
26     int k = blockIdx.x * blockDim.x + threadIdx.x;
27     if (k >= N) return;
28
29     float sum_real = 0.0f;
30     float sum_imag = 0.0f;
31
32     float angle_base = -2.0f * PI * k / N;
33
34     for (int tile = 0; tile < N; tile += BLOCK) {
35
36         int idx = tile + threadIdx.x;
37         if (idx < N) {
38             sh_real[threadIdx.x] = in_real[idx];
39             sh_imag[threadIdx.x] = in_imag[idx];
40         } else {
41             sh_real[threadIdx.x] = 0.0f;
42             sh_imag[threadIdx.x] = 0.0f;
43         }
44
45         __syncthreads();
46
47         int tileSize = min(BLOCK, N - tile);
48         for (int n = 0; n < tileSize; n++) {
49             float angle = angle_base * (tile + n);
50             float s, c;
51             sincosf(angle, &s, &c);
52
53             sum_real += sh_real[n] * c - sh_imag[n] * s;
54             sum_imag += sh_real[n] * s + sh_imag[n] * c;
55         }
56
57         __syncthreads();
58     }
59
60     out_real[k] = sum_real;
61     out_imag[k] = sum_imag;
62 }
63 int main(int argc, char** argv) {
64     int N = 4096;
65     if (argc >= 2) N = std::atoi(argv[1]);
66     size_t bytes = N * sizeof(float);
67
68     std::vector<float> h_real(N, 1.0f);
69     std::vector<float> h_imag(N, 0.0f);
70     std::vector<float> h_out_real(N), h_out_imag(N);
71
72     float *d_real, *d_imag, *d_out_real, *d_out_imag;
73     CHECK_CUDA(cudaMalloc(&d_real, bytes));
74     CHECK_CUDA(cudaMalloc(&d_imag, bytes));
75     CHECK_CUDA(cudaMalloc(&d_out_real, bytes));
76     CHECK_CUDA(cudaMalloc(&d_out_imag, bytes));
77
78     CHECK_CUDA(cudaMemcpy(d_real, h_real.data(), bytes,
79                           cudaMemcpyHostToDevice));
80     CHECK_CUDA(cudaMemcpy(d_imag, h_imag.data(), bytes,
81                           cudaMemcpyHostToDevice));

```



```

81     dim3 block(BLOCK);
82     dim3 grid((N + BLOCK - 1) / BLOCK);
83
84     cudaEvent_t start, stop;
85     cudaEventCreate(&start);
86     cudaEventCreate(&stop);
87
88     cudaEventRecord(start);
89     dft_cuda_shared<<<grid, block>>>(d_real, d_imag, d_out_real,
90     d_out_imag, N);
91     cudaEventRecord(stop);
92     cudaEventSynchronize(stop);
93
94     float ms;
95     cudaEventElapsedTime(&ms, start, stop);
96
97     CHECK_CUDA(cudaMemcpy(h_out_real.data(), d_out_real, bytes,
98     cudaMemcpyDeviceToHost));
99     CHECK_CUDA(cudaMemcpy(h_out_imag.data(), d_out_imag, bytes,
100    cudaMemcpyDeviceToHost));
101
102     double energy = 0.0;
103     for (int i = 0; i < N; i++)
104         energy += h_out_real[i] * h_out_real[i] + h_out_imag[i] *
105         h_out_imag[i];
106
107     std::cout << ms << std::endl;
108
109     cudaFree(d_real);
110     cudaFree(d_imag);
111     cudaFree(d_out_real);
112     cudaFree(d_out_imag);
113 }

```

### A.1.3 Înmulțirea Matricelor (gemm\_cuda.cu)

```

1  #include <iostream>
2  #include <vector>
3  #include <cuda_runtime.h>
4  #include < cublas_v2.h>
5
6  #define BLOCK 32
7
8  #define CUDA_CHECK(x) \
9      if((x) != cudaSuccess) { \
10         std::cout << "CUDA Error: " << cudaGetErrorString(x) << " at \
11         line " << __LINE__ << std::endl; exit(1); \
12     }
13
14 #define CUBLAS_CHECK(x) \
15     if((x) != CUBLAS_STATUS_SUCCESS) { \
16         std::cout << "cuBLAS Error at line " << __LINE__ << std::endl; \
17         exit(1); \
18     }
19
20 // Shared-memory GEMM kernel

```

```

19 __global__ void gemm_shared(const float* A, const float* B, float* C,
20 int N) {
21     __shared__ float As[BLOCK][BLOCK];
22     __shared__ float Bs[BLOCK][BLOCK];
23
24     int row = blockIdx.y * BLOCK + threadIdx.y;
25     int col = blockIdx.x * BLOCK + threadIdx.x;
26
27     float sum = 0.0f;
28
29     for(int t = 0; t < N; t += BLOCK) {
30         if(row < N && t + threadIdx.x < N)
31             As[threadIdx.y][threadIdx.x] = A[row*N + t + threadIdx.x];
32         else
33             As[threadIdx.y][threadIdx.x] = 0.0f;
34
35         if(col < N && t + threadIdx.y < N)
36             Bs[threadIdx.y][threadIdx.x] = B[(t + threadIdx.y)*N + col];
37         else
38             Bs[threadIdx.y][threadIdx.x] = 0.0f;
39
40         __syncthreads();
41
42         for(int k = 0; k < BLOCK; k++)
43             sum += As[threadIdx.y][k] * Bs[k][threadIdx.x];
44
45         __syncthreads();
46     }
47
48     if(row < N && col < N)
49         C[row*N + col] = sum;
50 }
51
52 int main(int argc, char** argv) {
53     int N = 4096;
54     if (argc >= 2) N = std::atoi(argv[1]);
55
56     size_t bytes = N*N*sizeof(float);
57     std::vector<float> hA(N*N, 1.0f), hB(N*N, 1.0f), hC(N*N, 0.0f);
58
59     float *dA, *dB, *dC;
60     CUDA_CHECK(cudaMalloc(&dA, bytes));
61     CUDA_CHECK(cudaMalloc(&dB, bytes));
62     CUDA_CHECK(cudaMalloc(&dC, bytes));
63
64     CUDA_CHECK(cudaMemcpy(dA, hA.data(), bytes, cudaMemcpyHostToDevice));
65     ;
66     CUDA_CHECK(cudaMemcpy(dB, hB.data(), bytes, cudaMemcpyHostToDevice));
67     ;
68
69     float ms = 0.0f;
70
71     if(N < 2048) {
72         dim3 threads(BLOCK, BLOCK);
73         dim3 blocks((N+BLOCK-1)/BLOCK, (N+BLOCK-1)/BLOCK);
74
75         cudaEvent_t start, stop;

```

```

74     cudaEventCreate(&start);
75     cudaEventCreate(&stop);
76     cudaEventRecord(start);
77
78     gemm_shared<<<blocks, threads>>>(dA,dB,dC,N);
79     cudaEventRecord(stop);
80     cudaEventSynchronize(stop);
81     cudaEventElapsedTime(&ms,start,stop);
82 } else {
83
84     cublasHandle_t handle;
85     CUBLAS_CHECK(cublasCreate(&handle));
86
87     float alpha = 1.0f;
88     float beta  = 0.0f;
89
90     cudaEvent_t start, stop;
91     cudaEventCreate(&start);
92     cudaEventCreate(&stop);
93     cudaEventRecord(start);
94
95
96     CUBLAS_CHECK(cublasSgemv(handle,
97         CUBLAS_OP_N, CUBLAS_OP_N,
98         N, N, N,
99         &alpha,
100        dB, N,
101        dA, N,
102        &beta,
103        dC, N));
104
105     cudaEventRecord(stop);
106     cudaEventSynchronize(stop);
107     cudaEventElapsedTime(&ms,start,stop);
108
109     cublasDestroy(handle);
110 }
111
112     CUDA_CHECK(cudaMemcpy(hC.data(), dC, bytes, cudaMemcpyDeviceToHost))
113 ;
114
115     double sum = 0.0;
116     for(float x : hC) sum += x;
117
118     std::cout << ms << std::endl;
119
120     cudaFree(dA);
121     cudaFree(dB);
122     cudaFree(dC);
123
124     return 0;
125 }

```

## A.1.4 Scripturi Bash de Automatizare

```

1 #!/bin/bash
2 set -e

```

```

3
4 OUT="linear_results.csv"
5 echo "algo,version,N,time_ms" > "$OUT"
6
7 DFT_BIN="./DFT"
8 GEMM_BIN="./GEMM"
9
10 DFT_SIZES=(512 1024 2048 4096 8192 16384 32768)
11 GEMM_SIZES=(512 1024 2048 4096 8192 16384)
12
13 # DFT
14 for N in "${DFT_SIZES[@]}"; do
15     echo "Running DFT N=$N"
16
17     # sequential
18     TIME=$(($DFT_BIN/dft_secv $N)
19     echo "DFT,seq,$N,$TIME" >> "$OUT"
20
21     # omp
22     TIME=$(($DFT_BIN/dft_cpu $N)
23     echo "DFT,omp,$N,$TIME" >> "$OUT"
24
25     # cuda
26     TIME=$(($DFT_BIN/dft_cuda $N)
27     echo "DFT,cuda,$N,$TIME" >> "$OUT"
28     echo "-----">> "$OUT"
29 done
30
31 # GEMM
32 for N in "${GEMM_SIZES[@]}"; do
33     echo "Running GEMM N=$N"
34
35     if [ $N -le 2048 ]; then
36         TIME=$(($GEMM_BIN/gemm_secv $N)
37         echo "GEMM,seq,$N,$TIME" >> "$OUT"
38     fi
39
40     TIME=$(($GEMM_BIN/gemm_cpu $N)
41     echo "GEMM,omp,$N,$TIME" >> "$OUT"
42
43     TIME=$(($GEMM_BIN/gemm_cuda $N)
44     echo "GEMM,cuda,$N,$TIME" >> "$OUT"
45     echo "-----">> "$OUT"
46 done

```

```

1 #!/bin/bash
2 set -e
3
4 OUT="conv2d_results.csv"
5 echo "algo,version,N,time_ms" > "$OUT"
6
7 # Adjust this path to where your executables are located
8 CONV_BIN="./Convolutie2D"
9
10 SIZES=(256 512 1024 2048 4096 8192 16384)
11
12 echo "Starting Conv2d Benchmarks..."
13

```

```

14 for N in "${SIZES[@]}"; do
15     echo "Running Conv2d N=$N"
16
17
18
19     TIME=$(($CONV_BIN/conv2d_secv $N)
20     echo "CONV,seq,$N,$TIME" >> "$OUT"
21
22
23     # 2. OpenMP
24     TIME=$(($CONV_BIN/conv2d_cpu $N)
25     echo "CONV,omp,$N,$TIME" >> "$OUT"
26
27     # 3. CUDA
28     TIME=$(($CONV_BIN/conv2d_cuda $N)
29     echo "CONV,cuda,$N,$TIME" >> "$OUT"
30
31     # Separator for readability (filtered out by python script later)
32     echo "-----" >> "$OUT"
33 done
34
35 echo "Done. Results saved to $OUT"

```