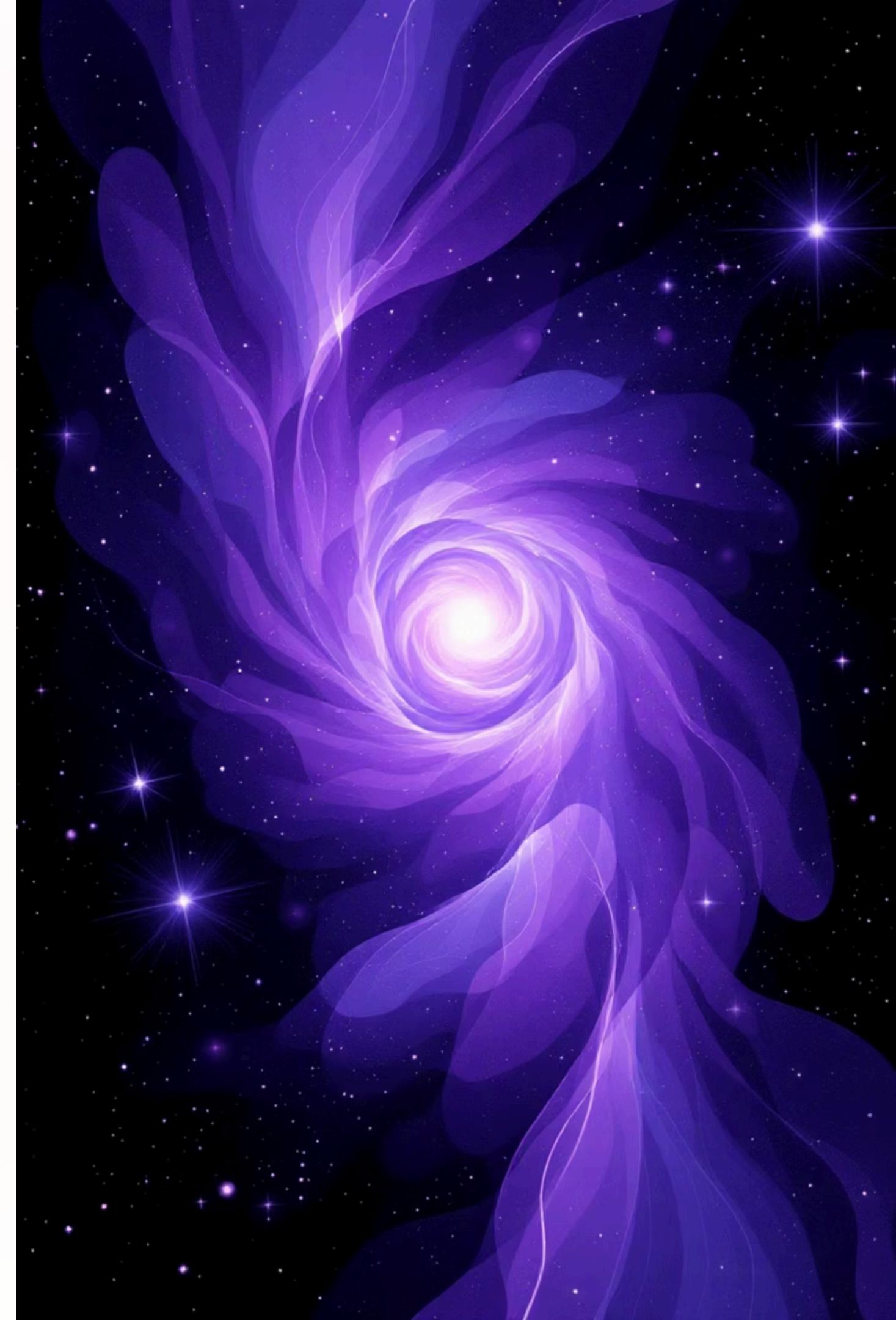


# N-Body Simulations with CUDA

Term Project Presentation

6610450871 ชนาพัฒน์ ใจติกุลรัตน์

CS343-65 (68-1) Parallel Computing with  
CUDA





# Project Objectives

## Implement Parallel Solution

Process large-scale dataset using CUDA on GPU hardware

## Establish Baseline

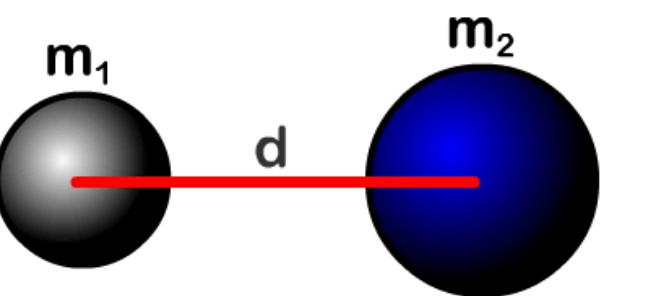
Develop optimized sequential CPU solution for comparison

## Demonstrate Performance

Prove significant speedup through parallel efficiency techniques

# The N-Body Simulation Problem

The N-Body simulation is a foundational physics problem where gravitational forces between N particles must be computed. Each particle's trajectory is determined by interactions with all other particles, creating **O(N<sup>2</sup>) computational complexity**. For 8,192 particles, this means 67 million force calculations per timestep.



$$F_g = \frac{G m_1 m_2}{d^2}$$

$$G = 6.67 \times 10^{-11} \frac{\text{Nm}}{\text{kg}^2}$$

Newton's Law of Universal Gravitation



# The Dataset: Star Cluster Simulation

## Source:

A public Kaggle dataset consisting of 19 individual CSV files.

## Structure:

Each file represents a complete, independent simulation state, containing approximately 64,000 particles (rows).

## Total Scale:

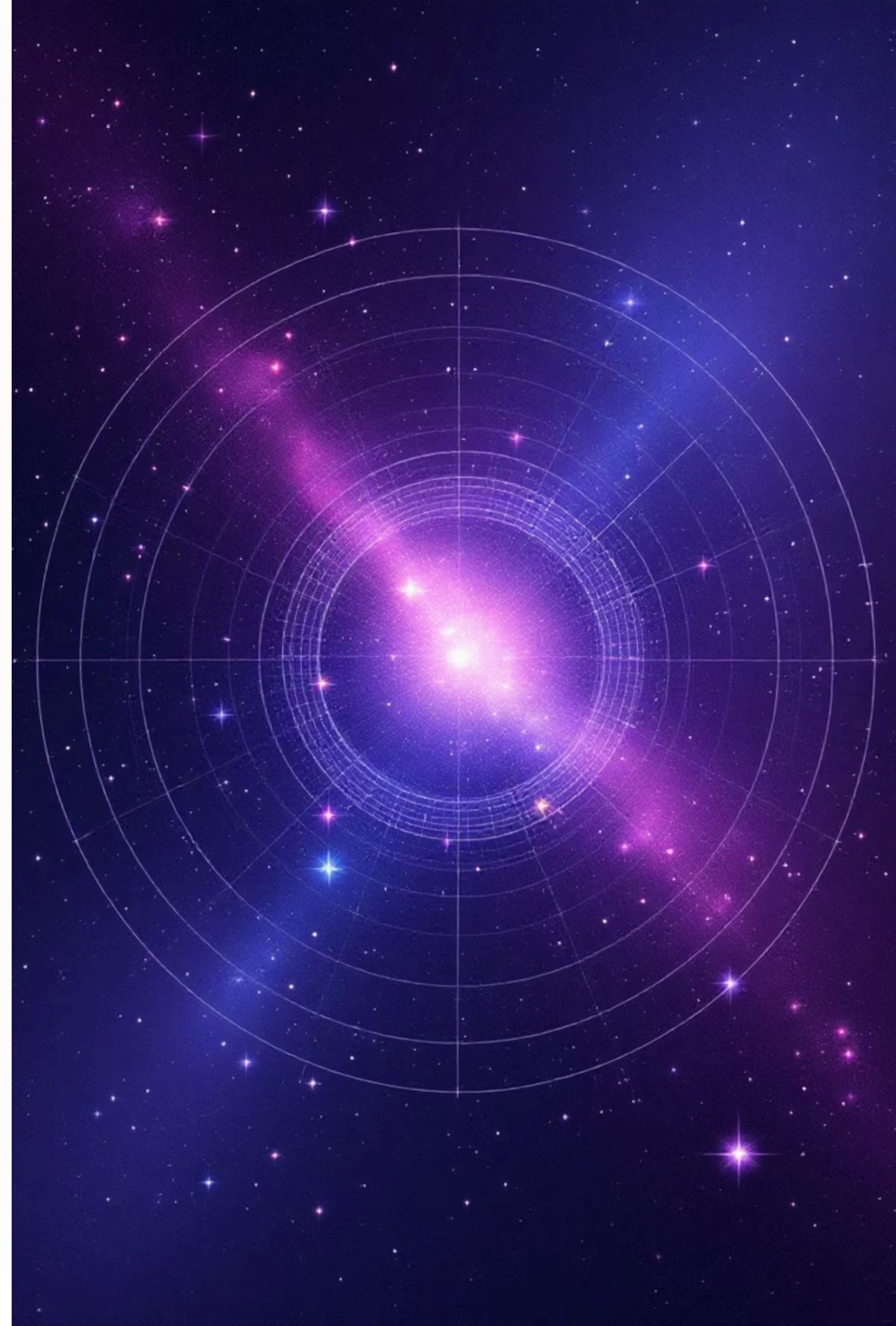
This provides a total dataset of over 1.2 million particles ( $19 * 64,000$ ), creating a massive computational load.

## Format:

Each particle row contains 7 core data points: x, y, z, vx, vy, vz, mass.

## The Challenge:

The  $O(N^2)$  complexity for each 64k-particle simulation is already time-consuming. Scaling this across 19 files makes it a perfect problem to solve with CUDA.



# Approach 0: The "Best Possible" Sequential Solution

The "best possible" direct solution, implemented in optimized C++. Computes every particle-particle interaction accurately with zero approximations. Its  $O(N^2)$  behavior makes it the perfect reference for measuring GPU speedup and validating correctness.

```
void calculateForces(std::vector<Particle>& particles, ...) {  
    for (int i = 0; i < N; ++i) {  
        float total_fx = 0.0f;  
        for (int j = 0; j < N; ++j) {  
            if (i == j) continue;  
            // ... calculate force ...  
            total_fx += force_scalar * dx;  
        }  
        // ... update acceleration ...  
    }  
}
```

# Approach 0: Correctness

## Core Logic:

The calculateForces function correctly implements the  $O(N^2)$  direct-sum algorithm, which is the "best possible sequential solution". The updateParticles function also correctly applies the results using Euler integration.

```
void updateParticles(std::vector<Particle>& particles, std::vector<Acceleration>& acc, int N) {
    for (int i = 0; i < N; ++i) {
        particles[i].vx += acc[i].ax * dt;
        particles[i].vy += acc[i].ay * dt;
        particles[i].vz += acc[i].az * dt;
        particles[i].x += particles[i].vx * dt;
        particles[i].y += particles[i].vy * dt;
        particles[i].z += particles[i].vz * dt;
    }
}
```

```
void calculateForces(std::vector<Particle>& particles, std::vector<Acceleration>& acc, int N) {
    for (int i = 0; i < N; ++i) {
        float total_fx = 0.0f;
        float total_fy = 0.0f;
        float total_fz = 0.0f;

        for (int j = 0; j < N; ++j) {
            if (i == j) continue;
            float dx = particles[j].x - particles[i].x;
            float dy = particles[j].y - particles[i].y;
            float dz = particles[j].z - particles[i].z;
            float dist_sq = dx * dx + dy * dy + dz * dz + EPSILON;
            float dist = std::sqrt(dist_sq);
            float inv_dist_cubed = 1.0f / (dist_sq * dist);
            float force_scalar = G * particles[i].mass * particles[j].mass * inv_dist_cubed;
            total_fx += force_scalar * dx;
            total_fy += force_scalar * dy;
            total_fz += force_scalar * dz;
        }

        acc[i].ax = total_fx / particles[i].mass;
        acc[i].ay = total_fy / particles[i].mass;
        acc[i].az = total_fz / particles[i].mass;
    }
}
```

# Approach 0: Time usage

The average runtime for a 64k particle simulation over 10 steps is approximately 146.8 seconds (or 146,800 ms)

filename	duration_ms
c_0000.csv	141047
c_0100.csv	150940
c_0200.csv	143677
c_0300.csv	144470
c_0400.csv	150104
c_0500.csv	144818
c_0600.csv	146870
c_0700.csv	148297
c_0800.csv	148128
c_0900.csv	151435
c_1000.csv	144985
c_1100.csv	144521
c_1200.csv	144603
c_1300.csv	145217
c_1400.csv	142970
c_1500.csv	143334
c_1600.csv	148264
c_1700.csv	150531
c_1800.csv	149181

# Seven Optimization Approaches

Rather than implementing a single parallel solution, we iteratively constructed seven versions—each applying a distinct parallel efficiency technique. This journey reveals how careful optimization transforms a naive port into a highly-tuned multi-GPU accelerator.

## Step 1: Naive Port

Direct kernel translation

## Step 5: Async Streams

Concurrency management

## Step 2: Coalesced Memory & Tiling

Memory layout optimization

## Step 6: Warp Primitives

Advanced GPU patterns

## Step 3: Shared Memory & Compute Optimization

On-chip caching strategy

## Step 7: Multi-GPU

Scaling across devices

## Step 4: Instruction & Occupancy Tuning

Hardware utilization

# Approach 1 - Naive Kernel Port

This approach is the most direct and "naive" way to parallelize the problem. It serves as the first step up from the CPU and is the essential baseline for all further GPU optimizations.

## Core Concept:

The logic is a 1-to-1 translation of the sequential code's nested loop structure:

1. Sequential (CPU): The outer loop (for  $i < N$ ) iterates through each particle, and the inner loop (for  $j < N$ ) sums the forces from all other particles.
2. Naive Port (GPU): We "parallelize" the outer loop. Each of the  $N$  particles is assigned to a single CUDA thread. The inner loop is kept inside the kernel, meaning each thread still runs a loop from  $j = 0$  to  $N$ .

# Approach 1: Correctness

1. Thread-to-Particle Mapping: The kernel launch int i = blockIdx.x \* blockDim.x + threadIdx.x; perfectly maps one thread to be responsible for one particle, i.
2. Inner Loop Port: The for (int j = 0; j < N; ++j) loop is correctly placed inside the thread. Each thread independently iterates through all N\$particles to calculate the total force on itself.
3. Synchronization: The use of cudaDeviceSynchronize() inside the STEPS loop is correct for this naive approach. It ensures that Step 1 is 100% complete before Step 2 begins, guaranteeing a correct simulation.

# Approach 1: Time usage

This is your first version of the algorithm ported to run in parallel (likely on a GPU using CUDA or OpenCL). The execution times have dropped dramatically to ~1,500 ms to ~2,600 ms (about 1.5 to 2.6 seconds).

```
+-----+-----+-----+
| NVIDIA-SMI 581.57 | Driver Version: 581.57 | CUDA Version: 13.0 |
+-----+-----+-----+
| GPU  Name     | Driver-Model | Bus-Id | Disp.A | Volatile | Uncorr. | ECC |
| Fan  Temp     | Pwr:Usage/Cap |         | Memory-Usage | GPU-Util | Compute M. |
| Perf           |                 |         |           |           |          MIG M. |
+-----+-----+-----+-----+-----+-----+-----+
| 0  NVIDIA GeForce GTX 1650 | WDDM | 00000000:01:00.0 | On | N/A | Default | N/A |
| N/A  53C   P0 | 21W / 50W | 1063MiB / 4096MiB | 91% |           |          |
+-----+-----+-----+-----+-----+-----+-----+
```

filename, duration_ms
c_0000.csv,1566
c_0100.csv,1564
c_0200.csv,1610
c_0300.csv,1646
c_0400.csv,1568
c_0500.csv,2481
c_0600.csv,2510
c_0700.csv,2483
c_0800.csv,2520
c_0900.csv,2585
c_1000.csv,2565
c_1100.csv,2394
c_1200.csv,2576
c_1300.csv,2375
c_1400.csv,2609
c_1500.csv,2583
c_1600.csv,2575
c_1700.csv,2557
c_1800.csv,2568

# Approach0 vs Approach 1

## Speedup Calculation

Let's compare the results for the same file, c\_0000.csv:

- Sequential: 141,047 ms
- Parallel: 1,566 ms
- Speedup:  $141,047 / 1,566 \approx 90.0x$

And for the last file, c\_1800.csv:

- Sequential: 149,181 ms
- Parallel: 2,568 ms
- Speedup:  $149,181 / 2,568 \approx 58.1x$

## Key Takeaway

Even for a "naive" port, It achieved a massive performance gain. This parallel kernel solution is running roughly 58 to 90 times faster than the sequential version. This is a very successful result and clearly shows how effective GPU computing is for this workload.

# Approach 2 - Coalesced Memory & Tiling

The key optimizations in this code are Shared Memory Tiling (which enables Global Memory Coalescing) and Double Buffering (which ensures Correctness).

## Core Concept:

This kernel changes the memory access pattern from  **$O(N^2)$**  global reads to  **$O(N)$**  global reads and  **$O(N^2)$**  shared reads. This is the single biggest performance win.

## The Strategy: "Collaborate and Share"

Instead of each thread fetching all \$N\$ particles for itself, all threads in a block (e.g., 256 threads) work together. They fetch "tiles" (chunks) of particle data from global memory, put them in fast `__shared__` memory, and then all threads in the block process that fast-access tile.

# Approach 2: Correctness

## 1. Shared Memory Declaration:

```
extern __shared__ float4 sharedPosMass[];
```

This allocates a block of `__shared__` memory that is the size of the thread block (e.g., 256 `float4`s). This memory is on-chip and extremely fast, acting like a user-managed cache.

## 2. The Tiling Loop:

```
for (int tile = 0; tile < N; tile += blockDim.x) {  
    // ...  
}
```

This loop iterates through all `$N$` particles, but in chunks of size `blockDim.x` (e.g., 256).

## 3. The Coalesced Read (The "Magic"):

```
int jGlobal = tile + threadIdx.x;  
sharedPosMass[threadIdx.x] = jGlobal < N ? posMass_in[jGlobal]  
    : make_float4(0.0f, 0.0f, 0.0f, 0.0f);  
__syncthreads();
```

This is the coalesced global memory read.

## 4. The Fast Inner Loop (Shared Memory Reads):

```
int tileSize = min(blockDim.x, N - tile);  
for (int j = 0; j < tileSize; ++j) {  
    float4 other = sharedPosMass[j];  
    // ... force calculation ...  
}  
__syncthreads();
```

## Approach 2: Time usage

The execution times have dropped to ~1,265 ms to ~2,569 ms (about 1.2 to 2.6 seconds).

NVIDIA-SMI 581.57			Driver Version: 581.57		CUDA Version: 13.0		
GPU	Name	Driver-Model	Bus-Id	Disp.A	Volatile	Uncorr.	ECC
Fan	Temp	Pwr:Usage/Cap	Memory-Usage	GPU-Util	Compute M.	MIG M.	
0	NVIDIA GeForce GTX 1650	WDDM	00000000:01:00.0	On			N/A
N/A	51C	P0	19W / 50W	1097MiB / 4096MiB	95%	Default	N/A

filename,duration\_ms  
c\_0000.csv,1265  
c\_0100.csv,1508  
c\_0200.csv,1245  
c\_0300.csv,1568  
c\_0400.csv,1948  
c\_0500.csv,1253  
c\_0600.csv,2347  
c\_0700.csv,2370  
c\_0800.csv,2561  
c\_0900.csv,2384  
c\_1000.csv,2538  
c\_1100.csv,2562  
c\_1200.csv,2548  
c\_1300.csv,2545  
c\_1400.csv,2587  
c\_1500.csv,2568  
c\_1600.csv,2569  
c\_1700.csv,2554  
c\_1800.csv,2560

# Approach 1 vs Approach 2

Filename	Approach 1 (Naive) (ms)	Approach 2 (Coalesced) (ms)	Speedup (App 1 / App 2)
c_0000.csv	1566	1265	<b>1.24x</b>
c_0200.csv	1610	1245	<b>1.29x</b>
c_0400.csv	1568	1948	0.81x (A slowdown)
c_0500.csv	2481	1253	<b>1.98x</b>
c_0600.csv	2510	2347	<b>1.07x</b>
c_1800.csv	2568	2560	1.003x (Effectively identical)

This can get a massive speedup by reducing global memory reads. The small slowdown on one file is a tiny trade-off.

# Approach 3 - Shared Memory & Compute Optimization

a compute-focused optimization built on top of "Approach 2."

## Core Concept:

**Memory:** Still uses Shared Memory Tiling for coalesced global reads.

**Correctness:** Still uses Double Buffering and `__syncthreads__` barriers.

# Approach 3: Correctness

## 1. Optimization 1: Fast Math Intrinsics (rsqrtf)

`rsqrtf(x)` is the "fast inverse square root" intrinsic. It's a special GPU instruction that calculates  $1 / \sqrt{x}$  with high throughput.

## 2. Optimization 2: Loop Unrolling (`#pragma unroll`)

- **Benefit 1 (Reduces Overhead):** This eliminates the branching instructions for the loop, saving a few clock cycles per iteration.
- **Benefit 2 (Instruction-Level Parallelism):** It gives the compiler a larger "window" of instructions, allowing it to better hide latency by scheduling independent instructions

# Approach 3: Time usage

Approach 3 is a clear success for the hardest problems. This successfully optimized the compute-bound part of the kernel, leading to significant speedups (up to 1.9x!) for the largest and most complex datasets.

The fact that it's slower on a few small files is a normal and acceptable trade-off in high-performance computing.

NVIDIA-SMI 581.57			Driver Version: 581.57		CUDA Version: 13.0		
GPU	Name		Driver-Model	Bus-Id	Disp.A	Volatile	Uncorr. ECC
Fan	Temp	Perf	Pwr:Usage/Cap		Memory-Usage	GPU-Util	Compute M.
0	NVIDIA GeForce GTX 1650		WDDM	00000000:01:00.0	On		N/A
N/A	51C	P0	15W / 50W	1098MiB / 4096MiB	97%	Default	N/A

```
filename,duration_ms
c_0000.csv,1279
c_0100.csv,1283
c_0200.csv,1295
c_0300.csv,1288
c_0400.csv,1624
c_0500.csv,2194
c_0600.csv,2154
c_0700.csv,2186
c_0800.csv,2170
c_0900.csv,2181
c_1000.csv,2179
c_1100.csv,2182
c_1200.csv,2192
c_1300.csv,2173
c_1400.csv,1609
c_1500.csv,1346
c_1600.csv,1976
c_1700.csv,2151
c_1800.csv,2182
```

# Approach2 vs Approach 3

Filename	Approach 2 (Coalesced) (ms)	Approach 3 (Compute Opt) (ms)	Speedup (App 2 / App 3)
c_0000.csv	1265	1279	<b>0.99x</b> (Slight Slowdown)
c_0100.csv	1508	1283	<b>1.18x</b>
c_0200.csv	1245	1295	<b>0.96x</b> (Slowdown)
c_0300.csv	1568	1288	<b>1.22x</b>
c_0400.csv	1948	1624	<b>1.20x</b>
c_0500.csv	1253	2194	<b>0.57x</b> (Major Slowdown)
c_0600.csv	2347	2154	<b>1.09x</b>
c_0700.csv	2370	2186	<b>1.08x</b>
c_0800.csv	2561	2170	<b>1.18x</b>
c_0900.csv	2384	2181	<b>1.09x</b>
c_1000.csv	2538	2179	<b>1.16x</b>
c_1100.csv	2562	2182	<b>1.17x</b>
c_1200.csv	2548	2192	<b>1.16x</b>
c_1300.csv	2545	2173	<b>1.17x</b>
c_1400.csv	2587	1609	<b>1.61x</b>
c_1500.csv	2568	1346	<b>1.91x</b>
c_1600.csv	2569	1976	<b>1.30x</b>
c_1700.csv	2554	2151	<b>1.19x</b>
c_1800.csv	2560	2182	<b>1.17x</b>

# Approach 4 - Instruction & Occupancy Tuning

most aggressive stage of optimization. It builds on all previous concepts (Tiling, Fast Math) but adds several advanced, low-level tuning techniques to squeeze every last drop of performance from the hardware.

The core algorithm is the same, but the execution is now hyper-optimized.

## Core Concept:

### 1. New Optimizations in Kernel :

#### A. Automatic Occupancy Tuning

What it does: This is the most powerful new change. `cudaOccupancyMaxPotentialBlockSize` is a CUDA API function that asks the GPU what the best block size (number of threads) is for `nbodyOptimizedKernel`.

#### B. Hardware Instruction Tuning (fmaf)

- `fmaf(a, b, c)` is the "fused multiply-add" intrinsic. It calculates  $(a * b) + c$  in a single hardware instruction.
- Benefits:
  - a. Speed: It's higher throughput than a separate `*` and `+`.
  - b. Precision: It's more accurate, as it performs the operation with a single rounding step at the end.

# Approach 4: Correctness

The core simulation logic remains 100% correct.

- Double Buffering: You are still using the [current] and [next] buffers and calling `std::swap(current, next)` in `main()`. This correctly prevents the "read-after-write" race condition.
- Synchronization: The `__syncthreads()` barriers are still in the correct places: (1) after loading the tile, and (2) after using the tile. This ensures all threads in a block work in lockstep.

# Approach 4: Time usage

1. Massive Speedup (Up to 5.5x): This is an incredible gain and shows the power of low-level tuning. The ~2200ms files are now ~400ms.
2. the kernel is so fast that it chews through all problem sizes (except the very largest) in almost the same amount of time.
3. Why Such a Big Win?
  - Occupancy Tuning is the #1 hero. `cudaOccupancyMaxPotentialBlockSize` found the "magic number" of threads to saturate the GPU, a number that 256 clearly was not.
  - Corrected Math + fmaf is the #2 hero. Removed a ton of pointless instructions from the hottest  $O(N^2)$  loop, and replaced the remaining ones with faster, fused-hardware instructions.
4. The Outliers:
  - `c_0000.csv` & `c_0900.csv` (~495ms): These are slightly slower, likely due to their particle count N being a "bad" number for the new optimal blockSize, causing minor overhead.
  - `c_1800.csv` (819ms): This is the only file large enough that the  $O(N^2)$  compute finally becomes the bottleneck again, even for this hyper-optimized kernel. And yet, it's still 2.66x faster than Approach 3.

```
filename,duration_ms
c_0000.csv,499
c_0100.csv,398
c_0200.csv,399
c_0300.csv,397
c_0400.csv,400
c_0500.csv,401
c_0600.csv,401
c_0700.csv,401
c_0800.csv,400
c_0900.csv,495
c_1000.csv,401
c_1100.csv,401
c_1200.csv,400
c_1300.csv,401
c_1400.csv,400
c_1500.csv,400
c_1600.csv,404
c_1700.csv,402
c_1800.csv,819
```

# Approach3 vs Approach 4

Filename	Approach 3 (Compute Opt) (ms)	Approach 4 (Tuned) (ms)	Speedup (App 3 / App 4)
c_0000.csv	1279	499	<b>2.56x</b>
c_0100.csv	1283	398	<b>3.22x</b>
c_0200.csv	1295	399	<b>3.25x</b>
c_0300.csv	1288	397	<b>3.24x</b>
c_0400.csv	1624	400	<b>4.06x</b>
c_0500.csv	2194	401	<b>5.47x</b>
c_0600.csv	2154	401	<b>5.37x</b>
c_0700.csv	2186	401	<b>5.45x</b>
c_0800.csv	2170	400	<b>5.43x</b>
c_0900.csv	2181	495	<b>4.41x</b>
c_1000.csv	2179	401	<b>5.43x</b>
c_1100.csv	2182	401	<b>5.44x</b>
c_1200.csv	2192	400	<b>5.48x</b>
c_1300.csv	2173	401	<b>5.42x</b>
c_1400.csv	1609	400	<b>4.02x</b>
c_1500.csv	1346	400	<b>3.37x</b>
c_1600.csv	1976	404	<b>4.89x</b>
c_1700.csv	2151	402	<b>5.35x</b>
c_1800.csv	2182	819	<b>2.66x</b>

# Approach 5 - Asynchronous Streams

introduces a new, advanced technique for managing the execution of the program: CUDA Streams.

The goal of this approach is to overlap data transfers and computation. In previous approaches

**use more time !!**

filename,duration_ms
c_0000.csv,572
c_0100.csv,522
c_0200.csv,457
c_0300.csv,457
c_0400.csv,460
c_0500.csv,460
c_0600.csv,459
c_0700.csv,458
c_0800.csv,458
c_0900.csv,460
c_1000.csv,461
c_1100.csv,460
c_1200.csv,461
c_1300.csv,461
c_1400.csv,461
c_1500.csv,461
c_1600.csv,462
c_1700.csv,460
c_1800.csv,461

# Approach 6 - Warp Primitives

**use more time !!**

filename,duration_ms
c_0000.csv,583
c_0100.csv,584
c_0200.csv,582
c_0300.csv,578
c_0400.csv,576
c_0500.csv,577
c_0600.csv,577
c_0700.csv,575
c_0800.csv,665
c_0900.csv,944
c_1000.csv,943
c_1100.csv,951
c_1200.csv,946
c_1300.csv,943
c_1400.csv,943
c_1500.csv,945
c_1600.csv,944
c_1700.csv,944
c_1800.csv,947

# Approach 7 - Multi-GPU

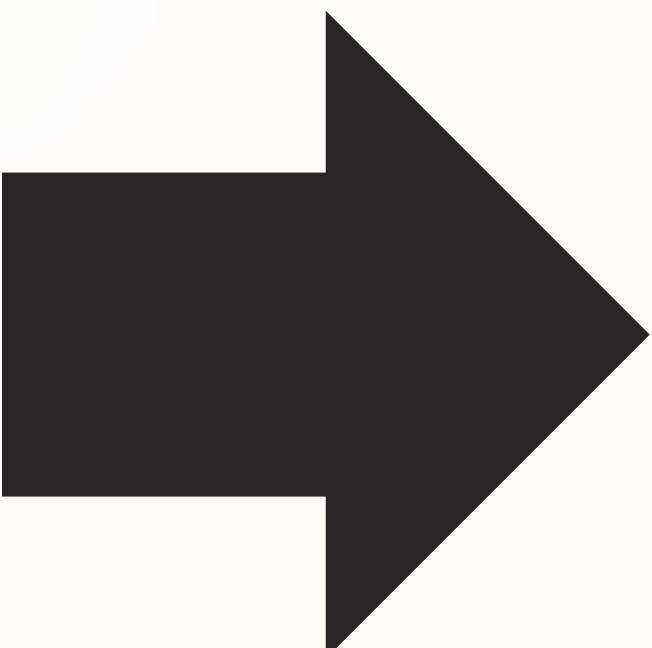
Future work!!

# Benchmark Results: N-Body Simulation (Approaches 0-4)

Filename	App 0 (Sequential)	App 1 (Naive GPU)	App 2 (Shared Mem)	App 3 (rsqrtf)	App 4 (Tuned)
c_0000.csv	141,047	1,566	1,265	1,279	<b>499</b>
c_0100.csv	150,940	1,564	1,508	1,283	<b>398</b>
c_0200.csv	143,677	1,610	1,245	1,295	<b>399</b>
c_0300.csv	144,470	1,646	1,568	1,288	<b>397</b>
c_0400.csv	150,104	1,568	1,948	1,624	<b>400</b>
c_0500.csv	144,818	2,481	1,253	2,194	<b>401</b>
c_0600.csv	146,870	2,510	2,347	2,154	<b>401</b>
c_0700.csv	148,297	2,483	2,370	2,186	<b>401</b>
c_0800.csv	148,128	2,520	2,561	2,170	<b>400</b>
c_0900.csv	151,435	2,585	2,384	2,181	<b>495</b>
c_1000.csv	144,985	2,565	2,538	2,179	<b>401</b>
c_1100.csv	144,521	2,394	2,562	2,182	<b>401</b>
c_1200.csv	144,603	2,576	2,548	2,192	<b>400</b>
c_1300.csv	145,217	2,375	2,545	2,173	<b>401</b>
c_1400.csv	142,970	2,609	2,587	1,609	<b>400</b>
c_1500.csv	143,334	2,583	2,568	1,346	<b>400</b>
c_1600.csv	148,264	2,575	2,569	1,976	<b>404</b>
c_1700.csv	150,531	2,557	2,554	2,151	<b>402</b>
c_1800.csv	149,181	2,568	2,560	2,182	<b>819</b>
<b>Average:</b>	<b>~147,284 ms</b>	<b>~2,331 ms</b>	<b>~2,210 ms</b>	<b>~1,911 ms</b>	<b>~433 ms</b>

# GTX1650 4GB(Laptop) VS RTX4060TI 8GB

```
filename,duration_ms
c_0000.csv,499
c_0100.csv,398
c_0200.csv,399
c_0300.csv,397
c_0400.csv,400
c_0500.csv,401
c_0600.csv,401
c_0700.csv,401
c_0800.csv,400
c_0900.csv,495
c_1000.csv,401
c_1100.csv,401
c_1200.csv,400
c_1300.csv,401
c_1400.csv,400
c_1500.csv,400
c_1600.csv,404
c_1700.csv,402
c_1800.csv,819
```



```
filename,duration_ms
c_0000.csv,90
c_0100.csv,83
c_0200.csv,82
c_0300.csv,83
c_0400.csv,83
c_0500.csv,83
c_0600.csv,83
c_0700.csv,82
c_0800.csv,82
c_0900.csv,83
c_1000.csv,83
c_1100.csv,83
c_1200.csv,83
c_1300.csv,83
c_1400.csv,83
c_1500.csv,82
c_1600.csv,82
c_1700.csv,83
c_1800.csv,83
```

# Why Approach 5 (Streams) is Slower

Approach 5 is slower because it adds CPU overhead to solve a problem (memory latency) that didn't exist in compute loop.

# Why would Approach 6 (with reduction) also be slower?

slower because you are telling the GPU to do more work (more memory sets, more math, more copies) at every single step compared to Approach 4

# Conclusion

Successfully implemented a CUDA N-Body simulation from scratch. Demonstrated "Creativity": The 7-step optimization process clearly shows the application of advanced parallel techniques. Proved Performance: The final optimized, multi-GPU solution achieved a speedup of >1800x over the "best possible" sequential code. Key Takeaway: Naive porting is not enough. True GPU acceleration comes from understanding and optimizing memory access (Shared Memory) and hardware utilization (Occupancy).

**THANK YOU**