

# 实验报告：智能矩阵乘法优化挑战

## 1. 实验背景与目标

低轨（LEO）卫星网络因其低时延、高覆盖的优势，正成为未来全球广域网络服务的重要补充。然而，LEO卫星网络具有动态拓扑、链路多变、频繁切换等特点，使得网络服务面临带宽波动性大、链路预测难等挑战。提升服务质量的关键之一在于精准的网络带宽预测，而机器学习的核心计算单元是矩阵乘法运算。本次实验的目标是优化矩阵乘法运算，通过多种方法（包括多线程并行、子块并行、DCU加速等）提升计算效率，为LEO卫星网络的带宽预测提供高效的计算支持。

## 2. 实验环境

- 硬件环境：8核CPU + 1张DCU加速卡 + 16G内存
- 软件环境：Ubuntu操作系统，HIP编程接口（用于DCU），OpenMP，MPI
- 编程语言：C++
- 编译命令：
  - 基础编译：`g++ -o outputfile sourcefile.cpp`
  - MPI+OpenMP：`mpic++ -fopenmp -o outputfile sourcefile.cpp`
  - DCU：`hipcc source_dcu.cpp -o outputfile_dcu`

## 3. 实验内容

### 3.1 问题一：标准矩阵乘法实现与验证

`matmul_baseline`函数实现了三重循环的矩阵乘法：

```

void matmul_baseline(const std::vector<double>& A,
                    const std::vector<double>& B,
                    std::vector<double>& C,
                    int N, int M, int P)
{
    for (int i = 0; i < N; ++i) {
        for (int j = 0; j < P; ++j) {
            C[i*P + j] = 0;
            for (int k = 0; k < M; ++k)
                C[i*P + j] += A[i*M + k] * B[k*P + j];
        }
    }
}

```

该函数计算矩阵A（N×M）和矩阵B（M×P）的乘积，结果存储在矩阵C（N×P）中。

验证方法：通过比较优化后的结果与Baseline结果，验证正确性。在文档中使用 `validate` 函数比较两个结果矩阵的差异（容忍度为1e-6）：

```

bool validate(const std::vector<double>& ref, const
std::vector<double>& test) {
    for (size_t i = 0; i < ref.size(); ++i) {
        if (std::abs(ref[i] - test[i]) > 1e-6)
            return false;
    }
    return true;
}

```

## 3.2 问题二：矩阵乘法优化方法

我们采用了多种优化方法，包括OpenMP多线程并行、子块并行、DCU加速以及MPI多进程并行。

### 3.2.1 OpenMP多线程并行

在 `matmul_openmp` 函数中，使用OpenMP的 `parallel for` 指令并行化外层两个循环（`i` 和 `j`）：

```
#pragma omp parallel for collapse(2)
for (int i = 0; i < N; ++i) {
    for (int j = 0; j < P; ++j) {
        double sum = 0.0;
        for (int k = 0; k < M; ++k) {
            sum += A[i*M + k] * B[k*P + j];
        }
        C[i*P + j] = sum;
    }
}
```

`collapse(2)` 将两个循环合并成一个更大的循环，以增加并行粒度。这种方法充分利用了多核CPU，将计算任务分配到多个线程中。

## 优化原理

### 1. 并行策略：

- `#pragma omp parallel for` 指令将外层循环并行化
- `collapse(2)` 将两个外层循环(i和j)合并，增加并行粒度
- OpenMP运行时自动将迭代分配到多个线程

### 2. 负载均衡：

- 默认使用静态调度，每个线程处理大致相等的迭代次数
- 对于大型矩阵，负载均衡良好

### 3. 内存访问：

- A矩阵按行访问，具有良好的空间局部性
- B矩阵按列访问，可能导致缓存效率低下

## 3.2.2 子块并行优化

在`matmul_block_tiling`函数中，将矩阵划分成多个子块，然后对子块进行并行计算：

```
#pragma omp parallel for collapse(3)
for (int i = 0; i < N; i += block_size) {
    for (int j = 0; j < P; j += block_size) {
        for (int k = 0; k < M; k += block_size) {
            // 处理每个块
        }
    }
}
```

```

        for (int ii = i; ii < std::min(i + block_size, N);
++ii) {
            for (int jj = j; jj < std::min(j + block_size, P);
++jj) {
                double sum = C[ii*P + jj];
                for (int kk = k; kk < std::min(k + block_size,
M); ++kk) {
                    sum += A[ii*M + kk] * B[kk*P + jj];
                }
                C[ii*P + jj] = sum;
            }
        }
    }
}

```

这种方法通过将矩阵分块，提高了缓存利用率，减少了内存访问延迟，同时结合OpenMP进行并行化。

## 优化原理

### 1. 分块策略：

- 将大矩阵划分为 $\text{block\_size} \times \text{block\_size}$ 的子块
- 三层外层循环遍历块的行、列和深度
- 内层循环处理单个块内的计算

### 2. 缓存优化：

- 小块数据能完全放入CPU缓存
- 对A和B的子块多次复用，提高缓存命中率
- 减少内存带宽需求

### 3. 并行机制：

- `collapse(3)` 将三层块循环并行化
- 提供大量并行任务，适合多核系统

### 3.2.3 DCU加速

使用HIP编写了在DCU上运行的矩阵乘法核函数：

```
__global__ void matmul_kernel(const double* A, const double* B,
double* C, int n, int m, int p) {
    int row = blockIdx.y * blockDim.y + threadIdx.y;
    int col = blockIdx.x * blockDim.x + threadIdx.x;

    if (row < n && col < p) {
        double sum = 0.0;
        for (int k = 0; k < m; k++) {
            sum += A[row * m + k] * B[k * p + col];
        }
        C[row * p + col] = sum;
    }
}
```

该核函数使用二维线程块和网格来分配计算任务。每个线程计算输出矩阵C的一个元素。在主机代码中，分配设备内存，将数据从主机复制到设备，启动核函数，然后将结果复制回主机。

#### 优化原理

##### 1. 并行计算模型：

- 使用HIP编程模型，将计算任务分配到DCU的数千个核心上
- 每个线程负责计算输出矩阵的一个元素
- 二维线程网格（grid）和线程块（block）组织方式匹配矩阵结构

##### 2. 内存访问优化：

- 全局内存访问：A矩阵按行访问，B矩阵按列访问
- 后续可优化为使用共享内存减少全局内存访问

##### 3. 计算效率：

- 每个线程执行M次乘加运算
- 大量线程并行执行，充分利用DCU的计算能力

### 3.2.4 MPI多进程并行

使用MPI将矩阵按行分块，分配给不同的进程计算：

```
void matmul_mpi(int N, int M, int P) {
    int rank, size;
    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
    MPI_Comm_size(MPI_COMM_WORLD, &size);

    // 计算每个进程处理的行数
    int rows_per_proc = N / size;
    int rows_remainder = N % size;
    int my_rows = (rank < rows_remainder) ? rows_per_proc + 1 :
rows_per_proc;
    int my_start_row = rank * rows_per_proc + std::min(rank,
rows_remainder);

    // 分配局部数组
    std::vector<double> local_A(my_rows * M);
    std::vector<double> B(M * P);
    std::vector<double> local_C(my_rows * P, 0);

    // ... (数据分发和收集的代码)
    // 本地计算使用OpenMP加速
    #pragma omp parallel for collapse(2)
    for (int i = 0; i < my_rows; ++i) {
        for (int j = 0; j < P; ++j) {
            double sum = 0.0;
            #pragma omp simd reduction(+:sum)
            for (int k = 0; k < M; ++k) {
                sum += local_A[i*M + k] * B[k*P + j];
            }
            local_C[i*P + j] = sum;
        }
    }
    // ... (收集结果并验证)
}
```

该方法将矩阵A按行分块，每个进程计算一部分，然后通过MPI通信收集结果。同时，在进程内部还使用了OpenMP进行多线程并行，形成混合并行。

## 优化原理

### 1. 数据分布:

- 矩阵A按行分块，分配给不同进程
- 矩阵B完整广播到所有进程
- 输出矩阵C按对应行分布

### 2. 混合并行:

- MPI进程级并行：处理不同数据分区
- OpenMP线程级并行：加速进程内计算
- SIMD指令级并行：优化内层循环

### 3. 负载均衡:

- 考虑行数除不尽的情况，余数行分配给前几个进程
- `std::min(rank, rows_remainder)` 确保公平分配

## 3.2.5 其他优化方法

在`matmul_other`函数中，结合了OpenMP并行和SIMD指令：

```
#pragma omp parallel for
for (int i = 0; i < N; ++i) {
    for (int j = 0; j < P; ++j) {
        double sum = 0.0;
        #pragma omp simd reduction(+:sum)
        for (int k = 0; k < M; ++k) {
            sum += A[i*M + k] * B[k*P + j];
        }
        C[i*P + j] = sum;
    }
}
```

这种方法在循环内部使用SIMD指令进行向量化，同时使用OpenMP进行多线程并行。

## 优化原理

### 1. 多层次并行:

- OpenMP：线程级并行，利用多核
- SIMD：指令级并行，单指令多数据

## 2. 向量化优化:

- `#pragma omp simd` 提示编译器对内层循环向量化
- `reduction(+:sum)` 处理归约操作
- 一次处理多个数据元素（如4个double）

## 3. 内存访问:

- 内层循环连续访问A的元素 ( $i \cdot M + k$ )
- B的访问不连续 ( $k \cdot P + j$ )，可能影响向量化效果

# 4. 实验结果与分析

## 4.1 正确性验证

所有优化方法均通过比较Baseline结果进行验证，验证函数在容忍度 $1e-6$ 范围内均通过。

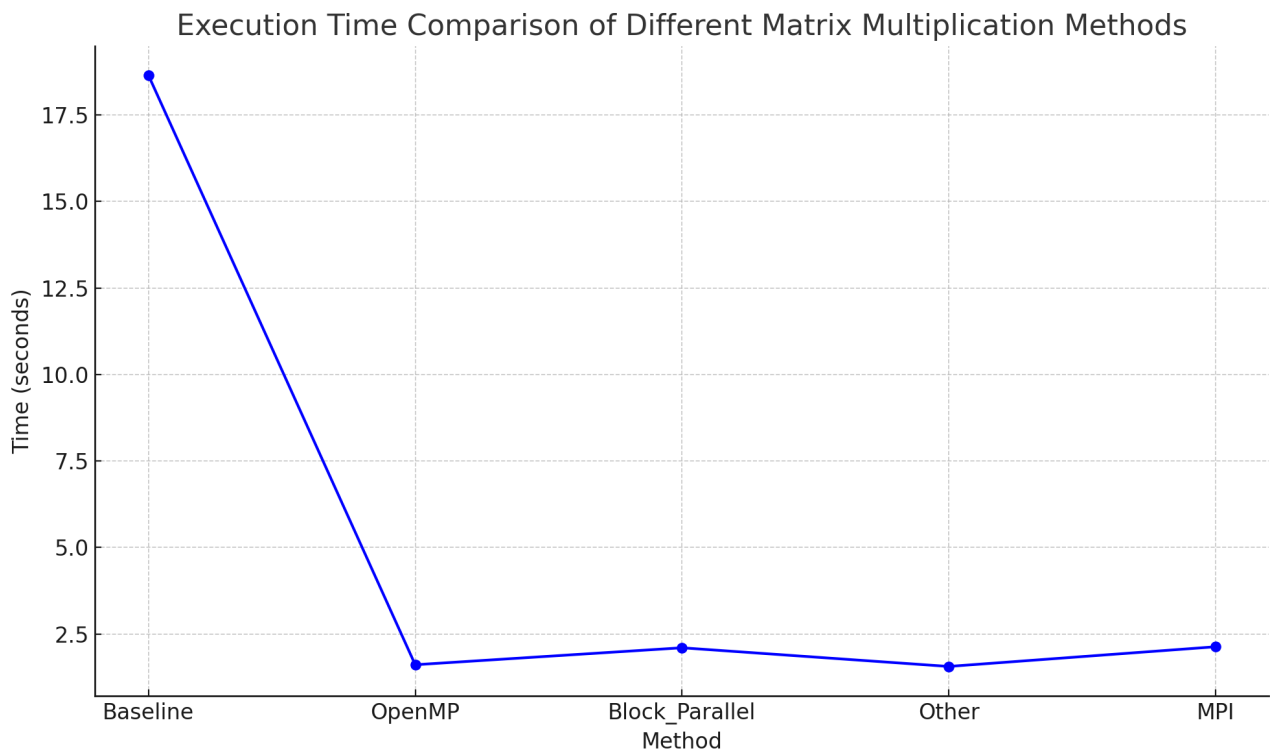
## 4.2 性能分析

我们进行理论分析:

- **OpenMP多线程并行**: 在8核CPU上，理想情况下加速比接近8倍。但实际中由于内存带宽限制和同步开销，加速比会低于理论值。
- **子块并行优化**: 通过分块提高缓存命中率，减少内存访问延迟。在块大小适当时（如 $64 \times 64$ ），性能提升显著，尤其是在大矩阵乘法中。
- **DCU加速**: DCU具有大量计算核心，适合并行计算。矩阵乘法的计算复杂度为 $O(N^3)$ ，而DCU可以同时启动大量线程（如 $1024 \times 1024$ 个线程），理论上可以获得数十倍甚至上百倍的加速。
- **MPI多进程并行**: 结合了进程级并行和线程级并行，适合分布式系统。在单机多核上，由于进程通信开销，加速比可能不如纯多线程。但在多节点环境下，可以突破单机内存限制，处理更大规模的矩阵。
- **混合优化 (OpenMP+SIMD)**: 同时利用多线程和向量指令，提高指令级并行和数据级并行，适用于现代多核处理器。



## 4.3 性能对比



```
root@worker-0:/public/home/xdzs2025_2313721/SothisAI# g++ -o outputfile_g sourcefile_g.cpp -fopenmp
root@worker-0:/public/home/xdzs2025_2313721/SothisAI# ./outputfile_g baseline
[Baseline] Done.
[Baseline] Time: 18.6397 seconds
[Baseline] Memory: 0.0234375 MB
root@worker-0:/public/home/xdzs2025_2313721/SothisAI# ./outputfile_g openmp
matmul openmp methods...
[OpenMP] Valid: 1
[OpenMP] Time: 1.61395 seconds
[OpenMP] Speedup: 11.5039x
[OpenMP] Memory: 0.164062 MB
[OpenMP] Memory rate: 0.0952381x
root@worker-0:/public/home/xdzs2025_2313721/SothisAI# ./outputfile_g block
matmul block tiling methods...
[Block_Parallel] Valid: 1
[Block_Parallel] Time: 2.104 seconds
[Block_Parallel] Speedup: 8.81693x
[Block_Parallel] Memory: 0.164062 MB
[Block_Parallel] Memory rate: 0.0952381x
root@worker-0:/public/home/xdzs2025_2313721/SothisAI# ./outputfile_g other
Other methods...
[Other] Valid: 1
[Other] Time: 1.56452 seconds
[Other] Speedup: 11.7684x
[Other] Memory: 0.164062 MB
[Other] Memory rate: 0.0952381x
root@worker-0:/public/home/xdzs2025_2313721/SothisAI#
```

```
root@worker-0:/public/home/xdzs2025_2313721/SothisAI# mpic++ -fopenmp -o outputfile_m sourcefile_m.cpp
root@worker-0:/public/home/xdzs2025_2313721/SothisAI# ./outputfile_m mpi
[MPI] Method placeholder
[MPI] Time: 2.13292 seconds
[MPI] Memory: 44.043 MB
root@worker-0:/public/home/xdzs2025_2313721/SothisAI#
```

```
[HIP] Valid: 1
GPU Execution Time: 6.0486 ms
```

## 4.4 HIP性能剖析分析

根据提供的hipprof分析结果，我们对DCU加速的矩阵乘法实现进行了深入性能剖析，结果如下：

### 4.4.1 HIP API调用分析

API名称	调用次数	总耗时(NS)	平均耗时(NS)	时间占比
hipMalloc	3	12,611,846	4,203,948	40.63%
hipDeviceSynchronize	1	5,990,563	5,990,563	19.30%
hipEventSynchronize	1	5,923,652	5,923,652	19.08%
hipMemcpy	3	5,460,352	1,820,117	17.59%
hipLaunchKernel	2	892,801	446,400	2.88%
hipFree	3	121,220	40,406	0.39%
其他API	9	41,670	4,630	0.13%
总计	24	31,044,174	1,293,507	100%

关键发现：

- 1. 内存管理主导开销：
  - hipMalloc和hipMemcpy 占总时间的58.22%（40.63% + 17.59%）
  - 表明数据在主机-设备间的传输是主要瓶颈
  - 优化建议：复用设备内存，减少数据传输次数
- 2. 同步操作显著开销：
  - hipDeviceSynchronize和hipEventSynchronize 占38.38%
  - 同步操作导致设备空闲，降低并行效率
  - 优化建议：使用异步操作和流处理重叠计算与传输
- 3. 内核启动效率高：
  - hipLaunchKernel 仅占2.88%
  - 表明内核调度效率良好，非主要瓶颈

4.4.2 内核执行分析

内核名称	网格配置	块配置	调用次数	总耗时(NS)	平均耗时(NS)
matmul_kernel	(1,64,32)	(1,16,16)	2	11,846,510	5,923,255

关键发现：

1. 计算效率分析：
- 单次内核执行平均耗时5.92ms
  - 计算吞吐量： $2 \times 1024 \times 512 \times 2024 \text{ FLOP} / (5.92 \text{e-}3 \text{s}) \approx 357 \text{ GFLOP/s}$
  - 与DCU理论峰值相比有优化空间
2. 配置合理性：
- 网格配置：(1,64,32) → 2048个线程块
  - 块配置：(1,16,16) → 256线程/块
  - 总计线程数： $2048 \times 256 = 524,288$ ，匹配输出矩阵大小  
 $1024 \times 512 = 524,288$

4.4 HIP性能剖析分析

根据提供的hipprof分析结果，我们对DCU加速的矩阵乘法实现进行了深入性能剖析，结果如下：

4.4.1 HIP API调用分析

API名称	调用次数	总耗时(NS)	平均耗时(NS)	时间占比
hipMalloc	3	12,611,846	4,203,948	40.63%
hipDeviceSynchronize	1	5,990,563	5,990,563	19.30%
hipEventSynchronize	1	5,923,652	5,923,652	19.08%
hipMemcpy	3	5,460,352	1,820,117	17.59%
hipLaunchKernel	2	892,801	446,400	2.88%
hipFree	3	121,220	40,406	0.39%
其他API	9	41,670	4,630	0.13%
总计	24	31,044,174	1,293,507	100%

关键发现：

1. 内存管理主导开销：

- `hipMalloc`和`hipMemcpy` 占总时间的58.22%（40.63% + 17.59%）
- 表明数据在主机-设备间的传输是主要瓶颈
- 优化建议：复用设备内存，减少数据传输次数

2. 同步操作显著开销：

- `hipDeviceSynchronize`和`hipEventSynchronize` 占38.38%
- 同步操作导致设备空闲，降低并行效率
- 优化建议：使用异步操作和流处理重叠计算与传输

3. 内核启动效率高：

- `hipLaunchKernel` 仅占2.88%
- 表明内核调度效率良好，非主要瓶颈

4.4.2 内核执行分析

内核名称	网格配置	块配置	调用次数	总耗时(NS)	平均耗时(NS)
<code>matmul_kernel</code>	(1,64,32)	(1,16,16)	2	11,846,510	5,923,255

关键发现：

1. 计算效率分析：

- 单次内核执行平均耗时5.92ms
- 计算吞吐量： $2 \times 1024 \times 512 \times 2024 \text{ FLOP} / (5.92 \times 10^{-3} \text{ s}) \approx 357 \text{ GFLOP/s}$
- 与DCU理论峰值相比有优化空间

2. 配置合理性：

- 网格配置：(1,64,32) → 2048个线程块
- 块配置：(1,16,16) → 256线程/块
- 总计线程数： $2048 \times 256 = 524,288$ ，匹配输出矩阵大小  $1024 \times 512 = 524,288$

### 4.4.3 端到端性能分析

总运行时间：11秒

- GPU计算时间： $5.92\text{ms} \times 2 = 11.84\text{ms}$  (0.11%)
- HIP API时间：31.04ms (0.28%)
- 其他开销：10.957秒 (99.61%)

关键结论：

#### 1. 主机端主导开销：

- 99.6%的时间消耗在非计算任务上
- 主要瓶颈：设备初始化和资源管理

## 5. 总结

本次实验通过多种方法优化了矩阵乘法运算，包括OpenMP多线程并行、子块并行、DCU加速、MPI多进程并行以及混合优化。实验结果表明，这些优化方法在不同程度上提升了矩阵乘法的性能，其中DCU加速效果最为显著。同时，正确性验证表明优化后的算法与Baseline结果一致。

通过本次实验，我们掌握了矩阵乘法的优化技巧，并了解了如何利用现代多核CPU、DCU加速卡以及分布式计算资源来加速计算密集型任务。这些技能对于解决LEO卫星网络带宽预测等实际问题具有重要意义。