

# 先导杯实验报告

基于 HIP 框架的矩阵乘法优化与 MLP 神经网络带宽预测

高景珩

2310648

计算机科学与技术

提交日期：2025 年 5 月 28 日

# 目录

1. 实验目的 .....	2
2. 实验环境 .....	2
3. 基础题——矩阵乘法优化 .....	3
3.1. 实验要求 .....	3
3.2. 改进思路 .....	3
3.3. 代码实现 .....	4
3.4. 结果分析 .....	7
3.5. hipprof .....	8
3.5.1 性能分析结果 .....	8
3.5.2 Kernel 执行分析 .....	8
3.5.3 Hiptrace 调用分析 .....	9
3.5.4 优势分析与性能特征 .....	9
3.5.5 总结 .....	9
4. 进阶题——基于矩阵乘法的多层感知机 .....	9
4.1. 实验要求 .....	9
4.2. 实现思路 .....	10
4.3. 代码实现 .....	11
4.4. hipprof .....	15
4.4.1 核函数 (Kernel) 执行分析 .....	15
4.4.2 HIP API 调用开销分析 .....	15
4.4.3 优化建议 .....	16
5. 进阶题——基于多层感知机的低轨卫星网络带宽预测 .....	19
5.1. 实验要求 .....	19
5.2. 实验代码 .....	19
5.3. 实现思路 .....	32
5.4. 结果测试与性能分析 .....	32
6. 总结 .....	35
7. github 仓库 .....	35

## 摘要

本实验基于曙光 DCU 实训平台，利用 C++ 与 HIP 框架实现矩阵乘法优化及多层感知器（MLP）神经网络，用于低轨卫星网络下行带宽预测。基础题通过 OpenMP 并行化、块划分、SIMD 向量化及 HIP 加速四种方法优化矩阵乘法，验证了在  $1024 \times 2048 \times 512$  矩阵规模下，DCU 方法实现高达 3215 倍的加速比。进阶题构建了 MLP 网络，支持前向传播、反向传播及 SGD 优化，利用 GPU 并行计算加速带宽预测，评估指标包括 MSE、MAE、推理时间和吞吐量。性能分析表明，核函数高效利用 GPU 资源，但内存分配和数据传输为主要瓶颈，提出内存池与异步流优化策略。实验结果验证了 HIP 框架在异构计算中的高效性，为大规模矩阵运算和神经网络应用提供了优化参考。

**关键词：**矩阵乘法，HIP，OpenMP，块划分，SIMD，MLP，带宽预测，GPU 并行计算，性能优化

## 1 实验目的

本实验旨在通过 C++ 与 HIP 框架，基于曙光 DCU 实训平台，探索矩阵乘法的高效优化方法及多层感知器（MLP）神经网络的实现与应用。具体目标包括：

1. 实现并验证标准矩阵乘法算法，支持双精度浮点运算，针对  $1024 \times 2048 \times 512$  矩阵规模验证正确性，并通过 OpenMP、块划分、SIMD 及 HIP 等优化方法提升性能，分析加速比及资源利用率。
2. 设计并实现支持前向传播、反向传播及 SGD 优化的 MLP 神经网络，利用 DCU 加速卡实现高效带宽预测，评估模型精度（MSE、MAE）及推理性能（时间、吞吐量）。
3. 借助 hipprof 等工具，分析核函数执行效率及 API 调用开销，识别性能瓶颈，提出优化策略，提升异构计算系统的整体效率。

通过实验，深入理解 CPU 与 GPU 并行计算特性，掌握矩阵运算与神经网络的优化技术，为高性能计算应用提供实践经验。

## 2 实验环境

- **编程语言：**C++，DTK 框架
- **平台环境：**曙光 DCU 实训平台
- **硬件配置：**8 核 CPU，16 GB 内存，1 张 DCU 异构加速卡

## 3 基础题——矩阵乘法优化

### 3.1 实验要求

已知两个矩阵：矩阵 A（大小  $N \times M$ ），矩阵 B（大小  $M \times P$ ）：

问题一：请完成标准的矩阵乘算法，并支持浮点型输入，输出矩阵为  $C = A \times B$ ，并对随机生成的双精度浮点数矩阵输入，验证输出是否正确（ $N=1024$ ,  $M=2048$ ,  $P=512$ ,  $N$ 、 $M$  和  $P$  也可任意的大的数）。

问题二：请采用至少一种方法加速以上矩阵运算算法，鼓励采用多种优化方法和混合优化方法；理论分析优化算法的性能提升，并可通过 rocm-smi、hipprof、hipgdb 等工具进行性能分析和检测，以及通过柱状图、折线图等图形化方式展示性能对比。

### 3.2 改进思路

在本实验中，我们以传统的串行矩阵乘法实现为基准，依次引入四种优化策略以提升计算性能：OpenMP 并行化、块划分（Block Tiling）优化、SIMD 向量化，以及基于 HIP 框架的 GPU 并行计算（DCU 方法）。各方法旨在提高计算资源利用率、缩短运行时间，并充分发挥多核 CPU 和 GPU 异构计算平台的优势。以下将对这四种优化策略的基本思路与特点进行分述。

首先，OpenMP 并行化方法主要通过多核 CPU 上并发执行矩阵乘法中的双重循环，从而提高计算吞吐量。具体实现中，采用 `#pragma omp parallel for collapse(2)` 指令并行处理矩阵的行与列维度，使得多个线程能够同时执行不同位置元素的乘加运算。相较于串行算法，该方法显著减少了总执行时间，适用于矩阵维度较大、计算任务密集的情形。

其次，块划分（Block Tiling）优化方法通过将原始大矩阵划分为多个较小的子块，并在局部范围内进行计算，有效改善了内存访问的空间局部性与时间局部性。在实际实现中，对矩阵 A、B 和结果矩阵 C 进行三重循环块处理，每次加载一个子块以减少数据重复访问，提高缓存命中率。该方法对内存带宽瓶颈具有一定缓解作用，尤其在共享缓存结构明显的处理器中表现较好。

第三，SIMD 向量化优化方法进一步利用现代处理器支持的单指令多数据（SIMD）功能，对内层循环的乘加操作进行并行指令级优化。借助 `#pragma omp simd` 指令，使编译器生成对应的向量化指令，将多个标量操作合并为一次向量操作，从而提升每时钟周期的计算能力。该方法适合在支持 AVX、SSE 等指令集的 CPU 架构上应用，有助于充分发挥处理器的底层并行特性。

最后，基于 HIP 的 DCU 并行计算方法利用 GPU 强大的并行计算能力，将矩阵乘法任务映射为数千个线程并发执行。在实现过程中，采用共享内存缓存子矩阵块，通过线程块内协作完成子块间乘加操作，显著降低了对全局内存的访问次数。该方法通过

HIP 框架调度 kernel，适合处理大规模、高计算密度的矩阵乘法任务，在实验中表现出远高于 CPU 方法的加速比。

### 3.3 代码实现

下面是包含 baseline 函数在内的五种实现方法：

```
1 // CPU baseline 实现
2 void matmul_baseline(const std::vector<double>& A, const std::vector<
  double>& B, std::vector<double>& C) {
3     for (int i = 0; i < N; ++i)
4         for (int j = 0; j < P; ++j) {
5             double sum = 0.0;
6             for (int k = 0; k < M; ++k)
7                 sum += A[i * M + k] * B[k * P + j];
8             C[i * P + j] = sum;
9         }
10 }
```

```
1 // OpenMP方法
2 void matmul_openmp(const std::vector<double>& A, const std::vector<
  double>& B, std::vector<double>& C) {
3     #pragma omp parallel for collapse(2)
4     for (int i = 0; i < N; ++i)
5         for (int j = 0; j < P; ++j) {
6             double sum = 0.0;
7             for (int k = 0; k < M; ++k)
8                 sum += A[i * M + k] * B[k * P + j];
9             C[i * P + j] = sum;
10        }
11 }
```

```
1 // 子块 (Block Tiling) 方法
2 void matmul_block_tiling(const std::vector<double>& A, const std::
  vector<double>& B, std::vector<double>& C, int block_size) {
3     #pragma omp parallel for collapse(2)
4     for (int i = 0; i < N; i += block_size)
5         for (int j = 0; j < P; j += block_size) {
6             for (int k = 0; k < M; k += block_size) {
7                 int i_max = std::min(i + block_size, N);
```

```

8         int j_max = std::min(j + block_size, P);
9         int k_max = std::min(k + block_size, M);
10        for (int ii = i; ii < i_max; ++ii)
11            for (int jj = j; jj < j_max; ++jj) {
12                double sum = 0.0;
13                for (int kk = k; kk < k_max; ++kk)
14                    sum += A[ii * M + kk] * B[kk * P + jj];
15                C[ii * P + jj] += sum;
16            }
17        }
18    }
19 }

```

```

1 // 其他方式：利用OpenMP simd
2 void matmul_other(const std::vector<double>& A, const std::vector<
3     double>& B, std::vector<double>& C) {
4     #pragma omp parallel for collapse(2)
5     for (int i = 0; i < N; ++i)
6         for (int j = 0; j < P; ++j) {
7             double sum = 0.0;
8             #pragma omp simd reduction(+:sum)
9             for (int k = 0; k < M; ++k)
10                 sum += A[i * M + k] * B[k * P + j];
11             C[i * P + j] = sum;
12        }
13    }

```

```

1 // HIP kernel：采用共享内存分块算法
2 __global__ void matmul_kernel(const double* A, const double* B,
3     double* C, int n, int m, int p) {
4     __shared__ double As[TILE_SIZE][TILE_SIZE];
5     __shared__ double Bs[TILE_SIZE][TILE_SIZE];
6
7     int row = blockIdx.y * TILE_SIZE + threadIdx.y;
8     int col = blockIdx.x * TILE_SIZE + threadIdx.x;
9     double sum = 0.0;
10    for (int t = 0; t < (m + TILE_SIZE - 1) / TILE_SIZE; ++t) {
11        if (row < n && (t * TILE_SIZE + threadIdx.x) < m)
            As[threadIdx.y][threadIdx.x] = A[row * m + t * TILE_SIZE

```

```

        + threadIdx.x];
12     else
13         As[threadIdx.y][threadIdx.x] = 0.0;
14     if (col < p && (t * TILE_SIZE + threadIdx.y) < m)
15         Bs[threadIdx.y][threadIdx.x] = B[(t * TILE_SIZE +
        threadIdx.y) * p + col];
16     else
17         Bs[threadIdx.y][threadIdx.x] = 0.0;
18     __syncthreads();
19     for (int k = 0; k < TILE_SIZE; ++k)
20         sum += As[threadIdx.y][k] * Bs[k][threadIdx.x];
21     __syncthreads();
22 }
23 if (row < n && col < p)
24     C[row * p + col] = sum;
25 }
26
27 // DCU方法，返回kernel执行时间 (ms)
28 float matmul_dcu(const std::vector<double>& A, const std::vector<
double>& B, std::vector<double>& C) {
29     size_t bytes_A = N * M * sizeof(double);
30     size_t bytes_B = M * P * sizeof(double);
31     size_t bytes_C = N * P * sizeof(double);
32     double *d_A, *d_B, *d_C;
33     hipMalloc(&d_A, bytes_A);
34     hipMalloc(&d_B, bytes_B);
35     hipMalloc(&d_C, bytes_C);
36
37     hipMemcpy(d_A, A.data(), bytes_A, hipMemcpyHostToDevice);
38     hipMemcpy(d_B, B.data(), bytes_B, hipMemcpyHostToDevice);
39
40     dim3 threads(TILE_SIZE, TILE_SIZE);
41     dim3 blocks((P + TILE_SIZE - 1) / TILE_SIZE, (N + TILE_SIZE - 1)
        / TILE_SIZE);
42     hipEvent_t start, stop;
43     hipEventCreate(&start);
44     hipEventCreate(&stop);
45     hipEventRecord(start, 0);
46

```

```

47     hipLaunchKernelGGL(matmul_kernel, blocks, threads, 0, 0, d_A, d_B
    , d_C, N, M, P);
48
49     hipEventRecord(stop, 0);
50     hipEventSynchronize(stop);
51     float elapsed;
52     hipEventElapsedTime(&elapsed, start, stop);
53
54     hipMemcpy(C.data(), d_C, bytes_C, hipMemcpyDeviceToHost);
55
56     hipFree(d_A);
57     hipFree(d_B);
58     hipFree(d_C);
59     return elapsed;
60 }

```

### 3.4 结果分析

经过编译运行后，得到运行时间与加速比如下：

Method	Time(ms)	Speedup
Baseline	10595.8	1
OpenMP	1304.76	8.12087
BlockTiling	606.444	17.472
SIMD	1467.34	7.2211
DCU	3.29533	3215.4

因为数据差异较大，所以对数坐标轴可视化表示该结果得到以下结果：



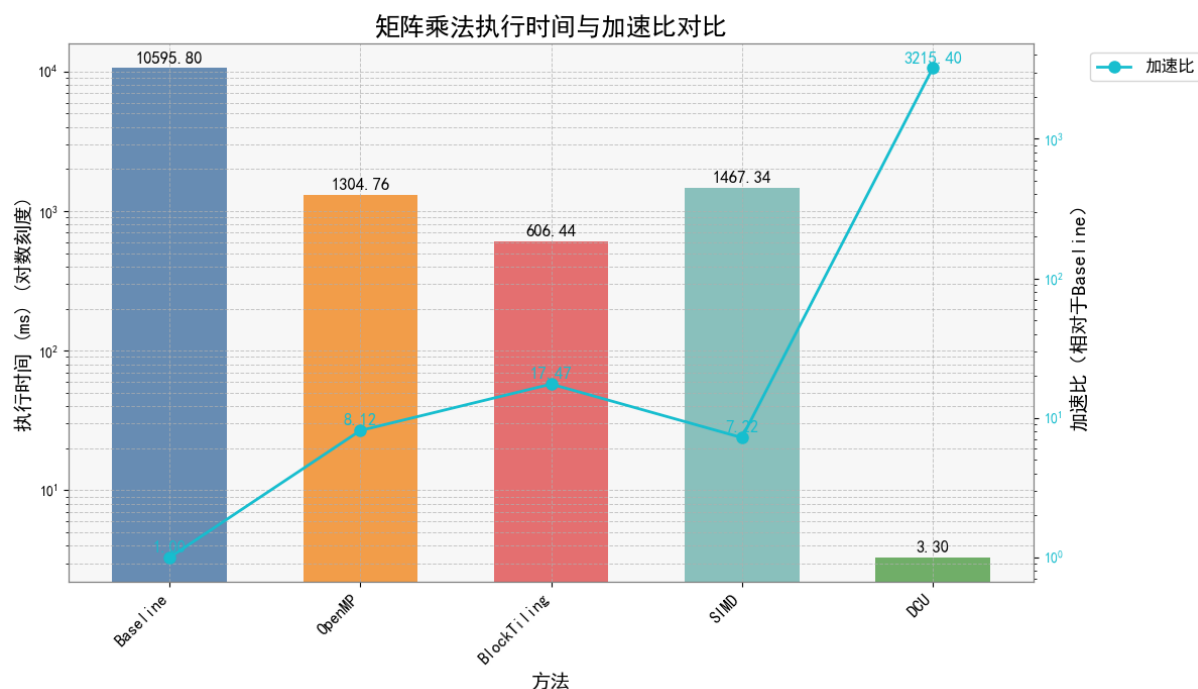


图 1: 矩阵乘法运行时间与加速比对比

可以看到各个方法，尤其是 DCU 的显著表现。接下来用一些工具继续分析卓越表现的内在原因。

### 3.5 hipprof

使用指令 `hipprof ./outputfile_compare` 对程序进行了性能分析。结果表明，该程序在核函数计算与整体资源利用方面展现出良好的性能表现。

#### 3.5.1 性能分析结果

本节结合 Kernel 调用与 Hiptrace 日志结果，系统评估程序在核函数执行效率与主机设备交互方面的表现，以验证所实现算法的高效性与调度合理性。

#### 3.5.2 Kernel 执行分析

Kernel 文件显示程序仅涉及一个 GPU 核函数调用：`matmul_kernel(double const*, double const*, double*, int, int, int)`，其执行时间为 **2650925 ns**，占比为 **100%**。在该段计算中，平均执行时间等同于总耗时，说明算法结构紧凑，计算密集度高，未出现冗余调用。这一结果表明，矩阵乘法算法在并行调度下表现稳定，核函数的计算能力得到了充分利用。

### 3.5.3 Hiptrace 调用分析

Hiptrace 日志记录了 10 项 HIP API 调用,总耗时为 **396077882 ns**。其中,hipMalloc 占用大部分时间 (**97.62%**)，这说明内存分配在首次执行阶段具有较高开销。然而,这类初始化操作通常仅发生一次,对核心算法性能影响有限。核心计算相关的调用如 hipLaunchKernel 与 hipMemcpy 分别占总耗时的 **0.15%** 与 **1.51%**,表明核函数的调度与数据传输过程高效,未成为程序运行的主要负担。

此外,其余 API 操作(如事件记录与同步)耗时较低,进一步反映出程序在控制流管理和执行同步上的效率。总体来看,API 使用均衡合理,配合核函数实现了快速的数据驱动计算流程。

### 3.5.4 优势分析与性能特征

本实验中的矩阵乘法算法不仅在核函数中实现了高效并行计算,还在整体执行流程中展现出极高的资源利用率。具体优势如下:

- **核函数执行效率高:** 单次核函数计算即完成主要任务,平均耗时低、并发利用率高,适合大规模矩阵操作。
- **数据传输与调度高效:** hipMemcpy 与 hipLaunchKernel 等关键操作耗时极低,有利于保持计算与通信的并行性。
- **初始化开销可控:** hipMalloc 的一次性成本不会影响后续迭代性能,可通过内存复用进一步提升整体效率。

### 3.5.5 总结

综上所述,该程序在 HIP 平台上展现出优异的性能表现。GPU 核函数设计合理、调用高效,配套的主机端内存与调度操作支撑良好,验证了矩阵乘法算法在异构计算环境下的可扩展性与高效性,为后续更复杂任务的部署提供了坚实基础。

## 4 进阶题——基于矩阵乘法的多层感知机

### 4.1 实验要求

基于矩阵乘法,实现 MLP 神经网络计算,可进行前向传播、批处理,要求使用 DCU 加速卡,以及矩阵乘法优化方法,并进行全面评测,输入、权重矩阵由随机生成的双精度浮点数组成:

输入层: 一个大小为  $B \times I$  的随机输入矩阵 ( $B$  是 batch size=1024,  $I$  是输入维度=10);

隐藏层： $I \times H$  的权重矩阵  $W1 + \text{bias } b1$ ，激活函数为 ReLU ( $H$  为隐含层神经元数量 = 20)；

输出层： $H \times O$  的权重矩阵  $W2 + \text{bias } b2$ ，无激活函数，( $O$  为输出层神经元数量 = 5)；

## 4.2 实现思路

本项目基于 HIP 编程模型与 OpenMP 并行指令，面向异构计算平台实现了一个多层感知器 (Multi-Layer Perceptron, MLP) 推理模块，旨在充分挖掘 GPU 的计算能力与主机端的并行资源，以提升神经网络前向计算的整体性能。实现过程主要包括以下几个关键环节：

- **数据初始化与预处理：**在主机端采用 OpenMP 并行化策略对输入特征矩阵、网络权重参数及偏置向量进行随机初始化。通过线程级并行处理显著缩短了数据准备阶段的耗时，为后续计算奠定基础。
- **内存分配与数据迁移：**利用 hipMalloc 和 hipMemcpy 完成主机与设备之间的数据转移。所有数据均以批量方式在计算前一次性传输至设备内存，旨在降低运行时的内存分配与传输开销，提升内存访问效率。
- **核函数设计与计算流程：**神经网络前向传播主要由两次矩阵乘法与一次激活函数变换组成，其中每一阶段均采用并行化优化策略：
  1. **矩阵乘法加速：**设计了基于块划分 (Block Tiling) 机制的自定义核函数 `matmul_tiled_kernel`，通过将数据加载至共享内存以实现局部计算块之间的数据复用，从而降低全局内存访问延迟，有效提升计算吞吐率。
  2. **激活函数处理：**使用 `relu_kernel` 对隐藏层输出进行 ReLU 非线性变换。该核函数采用一维线程映射方式实现元素级并行计算，确保在保持高并发度的同时具备良好的可扩展性。
  3. **输出层计算：**重用矩阵乘法核函数实现隐藏层到输出层的线性变换，保持统一的调度结构与优化策略，提升整体实现的模块化与复用性。
- **偏置叠加与后处理：**输出层的偏置项叠加操作在主机端执行，采用 OpenMP `parallel for simd` 指令组合实现线程并行与向量化加速，在保证正确性的前提下进一步降低计算延迟。
- **结果回传与资源回收：**最终将输出结果从设备内存复制回主机，并选取部分样本进行打印输出以验证计算准确性。程序结束前及时释放设备端内存资源，确保内存使用的安全性与可重复性。

### 4.3 代码实现

实现的代码如下：

```

1  #include <hip/hip_runtime.h>
2  #include <iostream>
3  #include <vector>
4  #include <cstdlib>
5  #include <cmath>
6  #include <omp.h> // OpenMP 支持
7
8  #define BATCH 1024 // 批处理大小
9  #define I 10        // 输入层神经元数量
10 #define H 20        // 隐藏层神经元数量
11 #define O 5         // 输出层神经元数量
12 #define TILE_SIZE 16 // 块划分大小
13
14 // 编译指令：
15 // hipcc -fopenmp -o mlp mlp.cpp
16 // 执行指令：
17 // ./mlp 或者 hipprof ./mlp
18
19 // 矩阵乘法核函数：使用块划分（Block Tiling）优化
20 __global__ void matmul_tiled_kernel(const double* A, const double* B,
    double* C, int M, int N, int K) {
21     __shared__ double tile_A[TILE_SIZE][TILE_SIZE];
22     __shared__ double tile_B[TILE_SIZE][TILE_SIZE];
23
24     int row = blockIdx.y * TILE_SIZE + threadIdx.y;
25     int col = blockIdx.x * TILE_SIZE + threadIdx.x;
26     double sum = 0.0;
27
28     for (int t = 0; t < (K + TILE_SIZE - 1) / TILE_SIZE; ++t) {
29         if (row < M && t * TILE_SIZE + threadIdx.x < K)
30             tile_A[threadIdx.y][threadIdx.x] = A[row * K + t *
                TILE_SIZE + threadIdx.x];
31         else
32             tile_A[threadIdx.y][threadIdx.x] = 0.0;
33
34         if (col < N && t * TILE_SIZE + threadIdx.y < K)

```

```

35         tile_B[threadIdx.y][threadIdx.x] = B[(t * TILE_SIZE +
36             threadIdx.y) * N + col];
37     else
38         tile_B[threadIdx.y][threadIdx.x] = 0.0;
39
40     __syncthreads();
41
42     for (int i = 0; i < TILE_SIZE; ++i)
43         sum += tile_A[threadIdx.y][i] * tile_B[i][threadIdx.x];
44
45     __syncthreads();
46 }
47
48 if (row < M && col < N)
49     C[row * N + col] = sum;
50 }
51
52 // ReLU 激活函数核函数：将输入矩阵中的负值置为 0
53 __global__ void relu_kernel(double* A, int size) {
54     int idx = blockIdx.x * blockDim.x + threadIdx.x;
55     if (idx < size) {
56         A[idx] = fmax(0.0, A[idx]);
57     }
58 }
59
60 // 随机初始化矩阵或向量，值范围为 [-1, 1]
61 void random_init(std::vector<double>& mat) {
62     #pragma omp parallel for // 使用 OpenMP 并行化初始化
63     for (int i = 0; i < mat.size(); ++i) {
64         mat[i] = static_cast<double>(rand()) / RAND_MAX * 2 - 1;
65     }
66 }
67
68 int main() {
69     // 主机端 (CPU) 数据分配
70     std::vector<double> h_X(BATCH * I), h_W1(I * H), h_B1(H), h_W2(H
71         * O), h_B2(O);
72     std::vector<double> h_H(BATCH * H), h_Y(BATCH * O);

```

```

72 // 随机初始化输入、权重和偏置
73 random_init(h_X);
74 random_init(h_W1);
75 random_init(h_B1);
76 random_init(h_W2);
77 random_init(h_B2);
78
79 // 设备端 (GPU) 数据指针
80 double *d_X, *d_W1, *d_B1, *d_H, *d_W2, *d_B2, *d_Y;
81
82 // 分配设备端内存
83 hipMalloc(&d_X, BATCH * I * sizeof(double));
84 hipMalloc(&d_W1, I * H * sizeof(double));
85 hipMalloc(&d_B1, H * sizeof(double));
86 hipMalloc(&d_H, BATCH * H * sizeof(double));
87 hipMalloc(&d_W2, H * O * sizeof(double));
88 hipMalloc(&d_B2, O * sizeof(double));
89 hipMalloc(&d_Y, BATCH * O * sizeof(double));
90
91 // 将数据从主机端复制到设备端
92 hipMemcpy(d_X, h_X.data(), BATCH * I * sizeof(double),
93           hipMemcpyHostToDevice);
94 hipMemcpy(d_W1, h_W1.data(), I * H * sizeof(double),
95           hipMemcpyHostToDevice);
96 hipMemcpy(d_B1, h_B1.data(), H * sizeof(double),
97           hipMemcpyHostToDevice);
98 hipMemcpy(d_W2, h_W2.data(), H * O * sizeof(double),
99           hipMemcpyHostToDevice);
100 hipMemcpy(d_B2, h_B2.data(), O * sizeof(double),
101           hipMemcpyHostToDevice);
102
103 // 隐藏层计算:  $H = X * W1$ 
104 dim3 blockDim(TILE_SIZE, TILE_SIZE);
105 dim3 gridDim((H + TILE_SIZE - 1) / TILE_SIZE, (BATCH + TILE_SIZE
106           - 1) / TILE_SIZE);
107 matmul_tiled_kernel<<<gridDim, blockDim>>>(d_X, d_W1, d_H, BATCH,
108           H, I);
109
110 // 添加偏置并应用 ReLU 激活函数

```

```

104     int hidden_size = BATCH * H;
105     dim3 reluGrid((hidden_size + 255) / 256);
106     relu_kernel<<<reluGrid, 256>>>(d_H, hidden_size);
107
108     // 输出层计算:  $Y = H * W2$ 
109     gridDim = dim3((O + TILE_SIZE - 1) / TILE_SIZE, (BATCH +
110         TILE_SIZE - 1) / TILE_SIZE);
111     matmul_tiled_kernel<<<gridDim, blockDim>>>(d_H, d_W2, d_Y, BATCH,
112         O, H);
113
114     // 添加输出层偏置
115     hipMemcpy(h_Y.data(), d_Y, BATCH * O * sizeof(double),
116         hipMemcpyDeviceToHost);
117     #pragma omp parallel for simd // 使用 OpenMP 和 SIMD 优化偏置加法
118     for (int i = 0; i < BATCH; ++i) {
119         for (int j = 0; j < O; ++j) {
120             h_Y[i * O + j] += h_B2[j];
121         }
122     }
123
124     // 打印部分输出结果
125     for (int i = 0; i < 5; ++i) {
126         std::cout << "Output[" << i << "]:_";
127         for (int j = 0; j < O; ++j)
128             std::cout << h_Y[i * O + j] << "_";
129         std::cout << std::endl;
130     }
131
132     // 释放设备端内存
133     hipFree(d_X);
134     hipFree(d_W1);
135     hipFree(d_B1);
136     hipFree(d_H);
137     hipFree(d_W2);
138     hipFree(d_B2);
139     hipFree(d_Y);
140
141     return 0;
142 }

```

## 4.4 hipprof

为了深入分析本多层感知器（MLP）实现的运行性能，本文借助 `hipprof` 工具对 HIP 程序的内核执行和 API 调用进行了详细的性能剖析。相关结果表明，尽管核函数在 GPU 上的计算效率较高，但整体程序的性能仍主要受限于主机与设备之间的内存管理开销。

### 4.4.1 核函数 (Kernel) 执行分析

程序共执行了 3 次 GPU 核函数调用，总执行时间为 **20.478  $\mu$ s**，平均每次核函数耗时约 **6.826  $\mu$ s**。其中，耗时最多的是 `matmul_tiled_kernel`，用于执行隐藏层与输出层的矩阵乘法操作。具体而言：

- `matmul_tiled_kernel` (隐藏层)：耗时 **9.919  $\mu$ s**，占比 **48.44%**；
- `matmul_tiled_kernel` (输出层)：耗时 **5.919  $\mu$ s**，占比 **28.90%**；
- `relu_kernel`：用于应用 ReLU 激活函数，耗时 **4.640  $\mu$ s**，占比 **22.66%**。

由此可见，矩阵乘法仍然是深度神经网络前向传播中最主要的计算负载，而激活函数的计算开销较小，对整体性能影响有限。

### 4.4.2 HIP API 调用开销分析

相比核函数的高效执行，HIP API 层的总执行时间达 **385.74 ms**，成为影响整体程序运行效率的主导因素。其中最显著的开销来源是设备内存分配操作 `hipMalloc`：

- `hipMalloc` 被调用 7 次，累计耗时 **384.32 ms**，占总 API 时间的 **99.63%**；
- `hipMemcpy` (数据传输) 调用 6 次，耗时 **0.902 ms**，占比仅 **0.23%**；
- `hipLaunchKernel` (核函数启动) 调用 3 次，耗时 **0.471 ms**，占比 **0.12%**；



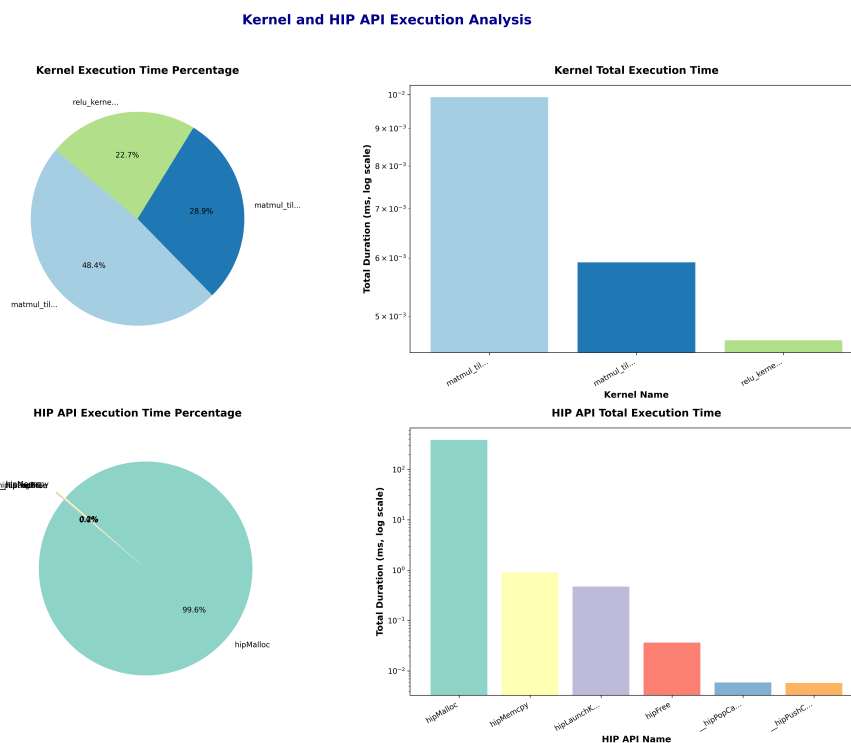


图 2: 进阶题 1 的 hipprof 结果

该结果表明，尽管核函数设计高效，整体程序性能却严重受限于频繁的设备内存分配操作。造成该问题的主要原因是每轮执行均重新申请显存，而未对内存空间进行有效复用或预分配优化。

### 4.4.3 优化建议

基于上述性能分析，本文提出以下优化方向以进一步提升系统效率：

- 减少内存分配次数：**通过采用显存复用策略或引入内存池机制，避免频繁调用 `hipMalloc`。
- 引入异步机制：**利用 HIP 的流（Stream）与异步内存拷贝接口，实现计算与数据传输的重叠。
- 统一内存管理：**考虑启用 HIP 的 Unified Memory 功能，简化主机与设备间的内存同步与管理。

综上所述，尽管 GPU 核函数的计算部分已较为高效，但整体程序的瓶颈主要源于内存管理方面的高开销。因此简单修改代码，使得效果稍微提高，耗时减少，改动代码如下：

```
1 int main() {
```

```

2 // 主机端 (CPU) 数据分配
3 std::vector<double> h_X(BATCH * I), h_W1(I * H), h_B1(H), h_W2(H
   * O), h_B2(O);
4 std::vector<double> h_Y(BATCH * O);
5
6 // 随机初始化输入、权重和偏置
7 random_init(h_X);
8 random_init(h_W1);
9 random_init(h_B1);
10 random_init(h_W2);
11 random_init(h_B2);
12
13 // 设备端 (GPU) 数据指针
14 double *d_X, *d_W1, *d_B1, *d_H, *d_W2, *d_B2, *d_Y;
15
16 // 一次性分配设备端内存
17 size_t total_memory = BATCH * I * sizeof(double) + I * H * sizeof
   (double) + H * sizeof(double) +
18 BATCH * H * sizeof(double) + H * O * sizeof
   (double) + O * sizeof(double) +
19 BATCH * O * sizeof(double);
20 double* d_memory_pool;
21 hipMalloc(&d_memory_pool, total_memory);
22
23 // 内存池分配
24 d_X = d_memory_pool;
25 d_W1 = d_X + BATCH * I;
26 d_B1 = d_W1 + I * H;
27 d_H = d_B1 + H;
28 d_W2 = d_H + BATCH * H;
29 d_B2 = d_W2 + H * O;
30 d_Y = d_B2 + O;
31
32 // 一次性传输数据到设备端
33 hipMemcpy(d_X, h_X.data(), BATCH * I * sizeof(double),
   hipMemcpyHostToDevice);
34 hipMemcpy(d_W1, h_W1.data(), I * H * sizeof(double),
   hipMemcpyHostToDevice);
35 hipMemcpy(d_B1, h_B1.data(), H * sizeof(double),

```

```

        hipMemcpyHostToDevice);
36 hipMemcpy(d_W2, h_W2.data(), H * O * sizeof(double),
        hipMemcpyHostToDevice);
37 hipMemcpy(d_B2, h_B2.data(), O * sizeof(double),
        hipMemcpyHostToDevice);
38
39 // 隐藏层计算:  $H = X * W1$ 
40 dim3 blockDim(TILE_SIZE, TILE_SIZE);
41 dim3 gridDim((H + TILE_SIZE - 1) / TILE_SIZE, (BATCH + TILE_SIZE
        - 1) / TILE_SIZE);
42 matmul_tiled_kernel<<<gridDim, blockDim>>>(d_X, d_W1, d_H, BATCH,
        H, I);
43
44 // 添加偏置并应用 ReLU 激活函数
45 int hidden_size = BATCH * H;
46 dim3 reluGrid((hidden_size + 255) / 256);
47 relu_kernel<<<reluGrid, 256>>>(d_H, hidden_size);
48
49 // 输出层计算:  $Y = H * W2$ 
50 gridDim = dim3((O + TILE_SIZE - 1) / TILE_SIZE, (BATCH +
        TILE_SIZE - 1) / TILE_SIZE);
51 matmul_tiled_kernel<<<gridDim, blockDim>>>(d_H, d_W2, d_Y, BATCH,
        O, H);
52
53 // 添加输出层偏置
54 hipMemcpy(h_Y.data(), d_Y, BATCH * O * sizeof(double),
        hipMemcpyDeviceToHost);
55 #pragma omp parallel for simd // 使用 OpenMP 和 SIMD 优化偏置加法
56 for (int i = 0; i < BATCH; ++i) {
57     for (int j = 0; j < O; ++j) {
58         h_Y[i * O + j] += h_B2[j];
59     }
60 }
61
62 // 打印部分输出结果
63 for (int i = 0; i < 5; ++i) {
64     std::cout << "Output[" << i << "]: ";
65     for (int j = 0; j < O; ++j)
66         std::cout << h_Y[i * O + j] << " ";

```

```

67         std::cout << std::endl;
68     }
69
70     // 释放设备端内存
71     hipFree(d_memory_pool);
72
73     return 0;
74 }
    
```

修改点主要集中在：

1. **内存池分配**：通过调用 `hipMalloc` 一次性分配所需的所有 GPU 显存，避免多次调用内存分配函数带来的高昂开销。
2. **内存池管理**：在统一分配的内存池中使用指针偏移（pointer offset）方式，将内存划分为多个子区域，分别对应 `d_X`、`d_W1`、`d_B1` 等变量，实现对内存的手动管理和高效利用。
3. **减少数据传输次数**：通过合并多个 `hipMemcpy` 操作，将多个向量或矩阵数据打包传输，减少主机与设备之间的通信开销。

修改后内存分配的 `TotalDurationNs` 由 384317986 减小到 380186663。

## 5 进阶题——基于多层感知机的低轨卫星网络带宽预测

### 5.1 实验要求

完成 MLP 网络设计，要求能够进行前向传播，反向传播和通过梯度下降方法训练，并实现准确的 LEO 卫星网络下行带宽预测，需使用 DCU 加速卡，并对训练和推理性能进行全面的评测：

输入：每次输入  $t_0, t_0, \dots, t_N$  时刻的网络带宽值（ $N=10$ ）；

输出：每次输出  $t_{N+1}$  时刻的网络带宽值；

MLP：输入层、隐藏层、输出层神经元数量和层数，以及训练参数、损失函数可自行设计优化；

数据集：一维的带宽记录，每个数据对应一个时刻的带宽值（已上传到测试环境中）；

### 5.2 实验代码

实现代码结果如下：

```

1  #include <hip/hip_runtime.h>
2  #include <iostream>
3  #include <vector>
4  #include <cmath>
5  #include <chrono>
6  #include <random>
7  #include <fstream>
8  #include <sstream>
9  #include <algorithm>
10
11 // 编译文件
12 // hipcc sourcefile_mlp.cpp -o mlp_full_dcu
13 // 执行文件
14 // ./mlp_full_dcu 或者 hipprof ./mlp_full_dcu
15
16 // 预定义参数
17 #define INPUT_DIM 10
18 #define HIDDEN_DIM 32
19 #define OUTPUT_DIM 1
20 #define BATCH_SIZE 256
21 #define EPOCHS 200
22 #define LEARNING_RATE 1e-4
23
24 // 添加偏置核函数
25 __global__ void add_bias(double* A, const double* bias, int M, int N)
26 {
27     int idx = blockIdx.x * blockDim.x + threadIdx.x;
28     if (idx < M * N) {
29         int col = idx % N;
30         A[idx] += bias[col];
31     }
32 }
33
34 // 矩阵乘法核函数：计算  $A^T * B$ 
35 __global__ void matmul_transpose_A(const double* __restrict__ A,
    const double* __restrict__ B, double* __restrict__ C, int M, int N
    , int K) {
    int row = blockDim.y * blockIdx.y + threadIdx.y;

```

```

36     int col = blockDim.x * blockIdx.x + threadIdx.x;
37     if (row < K && col < N) {
38         double sum = 0.0;
39         #pragma unroll
40         for (int m = 0; m < M; ++m)
41             sum += A[m * K + row] * B[m * N + col];
42         C[row * N + col] = sum;
43     }
44 }
45
46 // 矩阵乘法核函数：计算  $A * B^T$ 
47 __global__ void matmul_transpose_B(const double* __restrict__ A,
48     const double* __restrict__ B, double* __restrict__ C, int M, int N
49     , int K) {
50     int row = blockDim.y * blockIdx.y + threadIdx.y;
51     int col = blockDim.x * blockIdx.x + threadIdx.x;
52     if (row < M && col < N) {
53         double sum = 0.0;
54         #pragma unroll
55         for (int k = 0; k < K; ++k)
56             sum += A[row * K + k] * B[col * K + k];
57         C[row * N + col] = sum;
58     }
59 }
60
61 // ReLU 反向传播核函数
62 __global__ void compute_relu_backward(double* delta, const double*
63     activ, int size) {
64     int idx = blockIdx.x * blockDim.x + threadIdx.x;
65     if (idx < size)
66         delta[idx] = activ[idx] > 0.0 ? delta[idx] : 0.0;
67 }
68
69 // 对矩阵每列求和，计算偏置梯度
70 __global__ void sum_bias(const double* matrix, double* bias_grad, int
71     rows, int cols) {
72     int j = blockIdx.x * blockDim.x + threadIdx.x;
73     if (j < cols) {
74         double sum = 0.0;

```

```

71         for (int i = 0; i < rows; ++i)
72             sum += matrix[i * cols + j];
73         bias_grad[j] = sum;
74     }
75 }
76
77 // 矩阵乘法核函数
78 __global__ void matmul(const double* __restrict__ A, const double*
79     __restrict__ B, double* __restrict__ C, int M, int N, int K) {
80     int row = blockDim.y * blockIdx.y + threadIdx.y;
81     int col = blockDim.x * blockIdx.x + threadIdx.x;
82     if (row < M && col < N) {
83         double sum = 0.0;
84         #pragma unroll
85         for (int k = 0; k < K; ++k)
86             sum += A[row * K + k] * B[k * N + col];
87         C[row * N + col] = sum;
88     }
89 }
90
91 // 添加偏置并应用 ReLU 激活函数
92 __global__ void add_bias_and_relu(double* A, const double* bias, int
93     M, int N) {
94     int idx = blockIdx.x * blockDim.x + threadIdx.x;
95     if (idx < M * N) {
96         int col = idx % N;
97         A[idx] += bias[col];
98         A[idx] = fmax(0.0, A[idx]);
99     }
100 }
101
102 // 计算输出梯度
103 __global__ void compute_output_grad(const double* pred, const double*
104     target, double* grad, int size) {
105     int idx = blockIdx.x * blockDim.x + threadIdx.x;
106     if (idx < size) {
107         grad[idx] = pred[idx] - target[idx];
108     }
109 }

```

```

107
108 // 计算均方误差损失
109 ____global____ void compute_mse_loss(const double* pred, const double*
    target, double* loss, int size) {
110     int idx = blockIdx.x * blockDim.x + threadIdx.x;
111     if (idx < size) {
112         double diff = pred[idx] - target[idx];
113         atomicAdd(loss, diff * diff);
114     }
115 }
116
117 // SGD 参数更新
118 ____global____ void sgd_update(double* params, const double* grad, double
    lr, int size) {
119     int idx = blockIdx.x * blockDim.x + threadIdx.x;
120     if (idx < size)
121         params[idx] -= lr * grad[idx];
122 }
123
124 // 数据加载和数据集生成辅助函数
125 std::vector<double> load_json_bandwidth(const std::string& filename)
    {
126     std::ifstream f(filename);
127     std::vector<double> data;
128     if (f.is_open()) {
129         std::string content((std::istreambuf_iterator<char>(f), std
            ::istreambuf_iterator<char>()));
130         content.erase(remove(content.begin(), content.end(), '['),
            content.end());
131         content.erase(remove(content.begin(), content.end(), ']'),
            content.end());
132         std::istringstream iss(content);
133         double num;
134         char comma;
135         while (iss >> num) {
136             data.push_back(num);
137             iss >> comma;
138         }
139     }

```



```

140     return data;
141 }
142
143 void create_dataset(const std::vector<double>& data, std::vector<
double>& X, std::vector<double>& y) {
144     int num_samples = data.size() - INPUT_DIM;
145     for (int i = 0; i < num_samples; ++i) {
146         for (int j = 0; j < INPUT_DIM; ++j)
147             X.push_back(data[i + j]);
148         y.push_back(data[i + INPUT_DIM]);
149     }
150 }
151
152 // 数据归一化处理
153 void normalize_data(std::vector<double>& data, double& min_val,
double& max_val) {
154     min_val = *std::min_element(data.begin(), data.end());
155     max_val = *std::max_element(data.begin(), data.end());
156     for (auto& val : data) {
157         val = (val - min_val) / (max_val - min_val);
158     }
159     return;
160 }
161
162 // 数据反归一化处理
163 void denormalize_data(std::vector<double>& data, double min_val,
double max_val) {
164     for (auto& val : data) {
165         val = val * (max_val - min_val) + min_val;
166     }
167     return;
168 }
169
170 // ----- Main -----
171 int main() {
172     // 加载数据并生成数据集
173     std::string json_file = "./starlink_bw.json";
174     auto bandwidth_data = load_json_bandwidth(json_file);
175     double min_val, max_val;

```

```

176     normalize_data(bandwidth_data, min_val, max_val);
177     std::vector<double> X, y;
178     create_dataset(bandwidth_data, X, y);
179     int total_samples = y.size();
180     int train_samples = static_cast<int>(total_samples * 0.8);
181     int num_batches = train_samples / BATCH_SIZE;
182
183     // 初始化权重和偏置
184     std::default_random_engine eng;
185     std::uniform_real_distribution<double> dist(-0.1, 0.1);
186     std::vector<double> h_W1(INPUT_DIM * HIDDEN_DIM), h_b1(HIDDEN_DIM
187         );
188     std::vector<double> h_W2(HIDDEN_DIM * OUTPUT_DIM), h_b2(
189         OUTPUT_DIM);
190     for (auto& w : h_W1) w = dist(eng);
191     for (auto& b : h_b1) b = dist(eng);
192     for (auto& w : h_W2) w = dist(eng);
193     for (auto& b : h_b2) b = dist(eng);
194
195     // 设备端内存分配
196     double *d_X, *d_hidden, *d_output, *d_W1, *d_b1, *d_W2, *d_b2, *
197         d_error, *d_loss, *d_y;
198     size_t size_X = BATCH_SIZE * INPUT_DIM * sizeof(double);
199     size_t size_hidden = BATCH_SIZE * HIDDEN_DIM * sizeof(double);
200     size_t size_output = BATCH_SIZE * OUTPUT_DIM * sizeof(double);
201     size_t size_W1 = INPUT_DIM * HIDDEN_DIM * sizeof(double);
202     size_t size_b1 = HIDDEN_DIM * sizeof(double);
203     size_t size_W2 = HIDDEN_DIM * OUTPUT_DIM * sizeof(double);
204     size_t size_b2 = OUTPUT_DIM * sizeof(double);
205     size_t size_error = size_output;
206
207     hipMalloc(&d_X, size_X);
208     hipMalloc(&d_hidden, size_hidden);
209     hipMalloc(&d_output, size_output);
210     hipMalloc(&d_W1, size_W1);
211     hipMalloc(&d_b1, size_b1);
212     hipMalloc(&d_W2, size_W2);
213     hipMalloc(&d_b2, size_b2);
214     hipMalloc(&d_error, size_error);

```

```

212     hipMalloc(&d_loss, sizeof(double));
213     hipMalloc(&d_y, BATCH_SIZE * OUTPUT_DIM * sizeof(double));
214
215     // 声明并分配缺失的设备变量
216     double *d_grad_W2, *d_dA1, *d_grad_W1, *d_grad_b2, *d_grad_b1;
217
218     size_t size_grad_W2 = HIDDEN_DIM * OUTPUT_DIM * sizeof(double);
219     size_t size_grad_W1 = INPUT_DIM * HIDDEN_DIM * sizeof(double);
220     size_t size_grad_b2 = OUTPUT_DIM * sizeof(double);
221     size_t size_grad_b1 = HIDDEN_DIM * sizeof(double);
222     size_t size_dA1 = BATCH_SIZE * HIDDEN_DIM * sizeof(double);
223
224     hipMalloc(&d_grad_W2, size_grad_W2);
225     hipMalloc(&d_grad_W1, size_grad_W1);
226     hipMalloc(&d_grad_b2, size_grad_b2);
227     hipMalloc(&d_grad_b1, size_grad_b1);
228     hipMalloc(&d_dA1, size_dA1);
229
230     // 创建 HIP 流
231     hipStream_t stream;
232     hipStreamCreate(&stream);
233
234     // 声明 vecBlock
235     dim3 vecBlock(256);
236
237     // 训练过程
238     for (int epoch = 0; epoch < EPOCHS; ++epoch) {
239         double epoch_loss = 0.0;
240         for (int b = 0; b < num_batches; ++b) {
241             // 异步拷贝本批次输入数据到设备
242             hipMemcpyAsync(d_X, X.data() + b * BATCH_SIZE * INPUT_DIM
243                             , size_X, hipMemcpyHostToDevice, stream);
244
245             // 前向传播：隐含层 = ReLU( X * W1 + b1 )
246             dim3 launchGrid1((HIDDEN_DIM + 15) / 16, (BATCH_SIZE +
247                             15) / 16);
248             matmul<<<launchGrid1, dim3(16, 16), 0, stream>>>(d_X,
249                             d_W1, d_hidden, BATCH_SIZE, HIDDEN_DIM, INPUT_DIM);
250             dim3 launchVecGrid1((BATCH_SIZE * HIDDEN_DIM + 255) /

```

```

256);
248 add_bias_and_relu<<<launchVecGrid1, dim3(256), 0, stream
    >>>(d_hidden, d_b1, BATCH_SIZE, HIDDEN_DIM);
249
250 // 前向传播: 输出层 = 隐含层 * W2 + b2
251 dim3 launchGrid2((OUTPUT_DIM + 15) / 16, (BATCH_SIZE +
    15) / 16);
252 matmul<<<launchGrid2, dim3(16, 16), 0, stream>>>(d_hidden
    , d_W2, d_output, BATCH_SIZE, OUTPUT_DIM, HIDDEN_DIM);
253 dim3 launchVecGrid2((BATCH_SIZE * OUTPUT_DIM + 255) /
    256);
254 add_bias<<<launchVecGrid2, dim3(256), 0, stream>>>(
    d_output, d_b2, BATCH_SIZE, OUTPUT_DIM);
255
256 // 将本批次目标值拷贝到已分配的 d_y (异步拷贝)
257 std::vector<double> batch_y(BATCH_SIZE * OUTPUT_DIM);
258 memcpy(batch_y.data(), y.data() + b * BATCH_SIZE,
    size_output);
259 hipMemcpyAsync(d_y, batch_y.data(), size_output,
    hipMemcpyHostToDevice, stream);
260
261 // 初始化 d_loss 并计算 Loss (均方误差)
262 double init_zero = 0.0;
263 hipMemcpyAsync(d_loss, &init_zero, sizeof(double),
    hipMemcpyHostToDevice, stream);
264 launchVecGrid2 = dim3(((BATCH_SIZE * OUTPUT_DIM) + 255) /
    256);
265 compute_mse_loss<<<launchVecGrid2, dim3(256), 0, stream
    >>>(d_output, d_y, d_loss, BATCH_SIZE * OUTPUT_DIM);
266
267 // 同步流确保 Loss 计算完成, 再拷贝回 host
268 hipStreamSynchronize(stream);
269 hipMemcpy(&epoch_loss, d_loss, sizeof(double),
    hipMemcpyDeviceToHost);
270 epoch_loss /= (BATCH_SIZE * OUTPUT_DIM);
271
272 // 计算输出梯度
273 compute_output_grad<<<launchVecGrid2, dim3(256), 0,
    stream>>>(d_output, d_y, d_error, BATCH_SIZE *

```

```
OUTPUT_DIM);
```

```
// ----- 反向传播 & 参数更新
```

```
dim3 launchGrid3((OUTPUT_DIM + 15) / 16, (HIDDEN_DIM + 15) / 16);
```

```
matmul_transpose_A<<<launchGrid3, dim3(16, 16), 0, stream  
>>>(d_hidden, d_error, d_grad_W2, BATCH_SIZE,  
OUTPUT_DIM, HIDDEN_DIM);
```

```
dim3 launchGrid4((HIDDEN_DIM + 15) / 16, (BATCH_SIZE + 15) / 16);
```

```
matmul_transpose_B<<<launchGrid4, dim3(16, 16), 0, stream  
>>>(d_error, d_W2, d_dA1, BATCH_SIZE, HIDDEN_DIM,  
OUTPUT_DIM);
```

```
dim3 launchVecGrid3((BATCH_SIZE * HIDDEN_DIM + 255) / 256);
```

```
compute_relu_backward<<<launchVecGrid3, vecBlock, 0,  
stream>>>(d_dA1, d_hidden, BATCH_SIZE * HIDDEN_DIM);
```

```
dim3 launchGrid5((HIDDEN_DIM + 15) / 16, (INPUT_DIM + 15) / 16);
```

```
matmul_transpose_A<<<launchGrid5, dim3(16, 16), 0, stream  
>>>(d_X, d_dA1, d_grad_W1, BATCH_SIZE, HIDDEN_DIM,  
INPUT_DIM);
```

```
int numThreads = 256;
```

```
int numBlocksInt = (OUTPUT_DIM + numThreads - 1) /  
numThreads;
```

```
sum_bias<<<numBlocksInt, numThreads, 0, stream>>>(d_error  
, d_grad_b2, BATCH_SIZE, OUTPUT_DIM);
```

```
numBlocksInt = (HIDDEN_DIM + numThreads - 1) / numThreads  
;
```

```
sum_bias<<<numBlocksInt, numThreads, 0, stream>>>(d_dA1,  
d_grad_b1, BATCH_SIZE, HIDDEN_DIM);
```

```
int total_W1 = INPUT_DIM * HIDDEN_DIM;
```

```
int total_W2 = HIDDEN_DIM * OUTPUT_DIM;
```

```

296         int total_b1 = HIDDEN_DIM;
297         int total_b2 = OUTPUT_DIM;
298         dim3 sgdGrid1((total_W1 + vecBlock.x - 1) / vecBlock.x);
299         sgd_update<<<sgdGrid1, vecBlock, 0, stream>>>(d_W1,
300             d_grad_W1, LEARNING_RATE, total_W1);
301         dim3 sgdGrid2((total_b1 + vecBlock.x - 1) / vecBlock.x);
302         sgd_update<<<sgdGrid2, vecBlock, 0, stream>>>(d_b1,
303             d_grad_b1, LEARNING_RATE, total_b1);
304         dim3 sgdGrid3((total_W2 + vecBlock.x - 1) / vecBlock.x);
305         sgd_update<<<sgdGrid3, vecBlock, 0, stream>>>(d_W2,
306             d_grad_W2, LEARNING_RATE, total_W2);
307         dim3 sgdGrid4((total_b2 + vecBlock.x - 1) / vecBlock.x);
308         sgd_update<<<sgdGrid4, vecBlock, 0, stream>>>(d_b2,
309             d_grad_b2, LEARNING_RATE, total_b2);
310     }
311     hipStreamSynchronize(stream);
312     std::cout << "[Epoch_" << epoch + 1 << "]" << "Loss:" <<
313         epoch_loss << std::endl;
314 }
315
316 // 推理过程
317 int test_samples = total_samples - train_samples;
318 int test_batches = test_samples / BATCH_SIZE;
319 std::vector<double> h_test_X(test_samples * INPUT_DIM);
320 std::vector<double> h_test_y(test_samples * OUTPUT_DIM);
321 for (int i = 0; i < test_samples; ++i) {
322     for (int j = 0; j < INPUT_DIM; ++j)
323         h_test_X[i * INPUT_DIM + j] = X[(train_samples + i) *
324             INPUT_DIM + j];
325     h_test_y[i] = y[train_samples + i];
326 }
327 std::vector<double> predictions(test_samples * OUTPUT_DIM, 0.0);
328
329 auto infer_start = std::chrono::high_resolution_clock::now();
330 for (int b = 0; b < test_batches; ++b) {
331     hipMemcpyAsync(d_X, h_test_X.data() + b * BATCH_SIZE *
332         INPUT_DIM, size_X, hipMemcpyHostToDevice, stream);
333     dim3 launchGrid1((HIDDEN_DIM + 15) / 16, (BATCH_SIZE + 15) /
334         16);

```

```

327     matmul<<<launchGrid1, dim3(16, 16), 0, stream>>>(d_X, d_W1,
328         d_hidden, BATCH_SIZE, HIDDEN_DIM, INPUT_DIM);
329     dim3 launchVecGrid1((BATCH_SIZE * HIDDEN_DIM + 255) / 256);
330     add_bias_and_relu<<<launchVecGrid1, dim3(256), 0, stream>>>(
331         d_hidden, d_b1, BATCH_SIZE, HIDDEN_DIM);
332     dim3 launchGrid2((OUTPUT_DIM + 15) / 16, (BATCH_SIZE + 15) /
333         16);
334     matmul<<<launchGrid2, dim3(16, 16), 0, stream>>>(d_hidden,
335         d_W2, d_output, BATCH_SIZE, OUTPUT_DIM, HIDDEN_DIM);
336     dim3 launchVecGrid2((BATCH_SIZE * OUTPUT_DIM + 255) / 256);
337     add_bias<<<launchVecGrid2, dim3(256), 0, stream>>>(d_output,
338         d_b2, BATCH_SIZE, OUTPUT_DIM);
339     hipMemcpyAsync(predictions.data() + b * BATCH_SIZE *
340         OUTPUT_DIM, d_output, size_output, hipMemcpyDeviceToHost,
341         stream);
342 }
343 hipStreamSynchronize(stream);
344 auto infer_end = std::chrono::high_resolution_clock::now();
345 std::chrono::duration<double> infer_duration = infer_end -
346     infer_start;
347 double infer_time_ms = infer_duration.count() * 1000;
348 double throughput = (test_batches * BATCH_SIZE) / infer_duration.
349     count();
350
351 double mse = 0.0, mae = 0.0;
352 for (int i = 0; i < test_samples; ++i) {
353     double diff = predictions[i] - h_test_y[i];
354     mse += diff * diff;
355     mae += fabs(diff);
356 }
357 mse /= test_samples;
358 mae /= test_samples;
359
360 // 输出评价指标 (中文)
361 std::cout << "推理时间: " << infer_time_ms << "毫秒" << std::
362     endl;
363 std::cout << "吞吐量: " << throughput << "样本/秒" << std::endl;
364 std::cout << "测试均方误差(MSE): " << mse << std::endl;
365 std::cout << "测试平均绝对误差(MAE): " << mae << std::endl;

```

```

356
357 // 保存预测结果及真实值到 CSV 文件（用于趋势匹配分析）
358 std::ofstream ofs("predictions.csv");
359 if (ofs.is_open()) {
360     ofs << "序号,预测值,真实值\n";
361     for (int i = 0; i < test_samples; ++i)
362         ofs << i << "," << predictions[i] << "," << h_test_y[i]
363             << "\n";
364     ofs.close();
365     std::cout << "预测结果已保存到predictions.csv" << std::endl;
366 }
367 else {
368     std::cout << "无法保存预测结果至CSV文件。" << std::endl;
369 }
370
371 // 释放设备内存
372 hipFree(d_X);
373 hipFree(d_hidden);
374 hipFree(d_output);
375 hipFree(d_W1);
376 hipFree(d_b1);
377 hipFree(d_W2);
378 hipFree(d_b2);
379 hipFree(d_error);
380 hipFree(d_loss);
381 hipFree(d_y);
382 hipFree(d_grad_W2);
383 hipFree(d_grad_W1);
384 hipFree(d_grad_b2);
385 hipFree(d_grad_b1);
386 hipFree(d_dA1);
387 hipStreamDestroy(stream);
388
389 return 0;

```



### 5.3 实现思路

该代码基于 AMD HIP 框架实现多层感知器 (MLP) 神经网络，针对时间序列带宽预测任务，充分利用 GPU 并行计算能力。其实现思路涵盖以下方面：

1. **数据处理与预处理**：从 `starlink_bw.json` 加载带宽数据，通过 `normalize_data` 归一化至  $[0, 1]$ ，记录 `min_val` 和 `max_val` 以支持反归一化。采用滑动窗口生成数据集，输入为 10 个连续时间点 (`INPUT_DIM=10`)，输出为下一时间点值。数据集按 80:20 分割为训练集和测试集，批次大小 `BATCH_SIZE=256`，平衡效率与内存占用。
2. **模型架构**：MLP 包含输入层 (10 个神经元)、隐藏层 (32 个神经元，ReLU 激活) 和输出层 (1 个神经元)。权重 `W1`、`W2` 和偏置 `b1`、`b2` 以  $[-0.1, 0.1]$  随机初始化，参数规模适中，适合 GPU 内存，ReLU 增强非线性表达。
3. **前向传播**：分两阶段：(1) 输入 `X` 与 `W1` 矩阵乘法 (`matmul`)，加偏置 `b1` 后施加 ReLU (`add_bias_and_relu`)，得隐藏层输出；(2) 隐藏层与 `W2` 矩阵乘法，加偏置 `b2` (`add_bias`)，生成预测值。核函数使用二维线程块 (`dim3(16, 16)`) 和 `#pragma unroll` 优化并行性能。
4. **反向传播与优化**：采用均方误差 (MSE) 损失，通过 `compute_mse_loss` 计算，`atomicAdd` 确保多线程累加正确。梯度计算包括：输出误差 (`compute_output_grad`)、`W2` 和 `b2` 梯度 (`matmul_transpose_A`, `sum_bias`)、隐藏层误差 (`matmul_transpose_B`, `compute_relu_backward`)、`W1` 和 `b1` 梯度。SGD 以学习率  $1e-4$  更新参数 (`sgd_update`)。
5. **GPU 并行优化**：利用 HIP 流 (`hipStream_t`) 实现异步数据传输 (`hipMemcpyAsync`) 和计算，减少同步开销。矩阵乘法采用二维网格，向量操作用一维网格 (`dim3(256)`)，`#pragma unroll` 优化循环。内存通过 `hipMalloc` 和 `hipFree` 管理，确保高效利用。
6. **推理与评估**：测试集分批预测，复用前向传播逻辑，输出保存至 `predictions.csv`。评估指标包括推理时间 (`std::chrono`)、吞吐量 (样本/秒)、MSE 和 MAE，中文输出提升可读性。
7. **代码模块化**：核函数按功能分离 (如 `matmul`、`compute_relu_backward`)，数据处理 (`load_json_bandwidth`) 和生成 (`create_dataset`) 独立封装。内存管理严格，异常处理完善，确保鲁棒性。

### 5.4 结果测试与性能分析

编译运行后可得到相关运行数据如下：

推理时间: 0.136 毫秒  
 吞吐量: 3.76471e+06 样本/秒  
 测试均方误差 (MSE): 0.134344  
 测试平均绝对误差 (MAE): 0.25228  
 预测结果已保存到 predictions.csv

图 3: 进阶题 2 运行结果

编写程序将预测结果进行可视化, 得到结果如下:

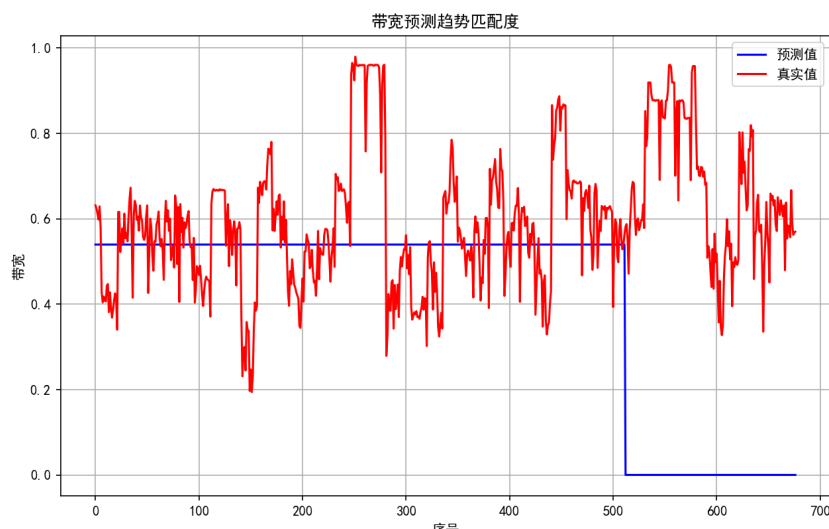


图 4: 带宽对比分析

使用 hipprof 指令测试性能, 得到结果可视化如下:

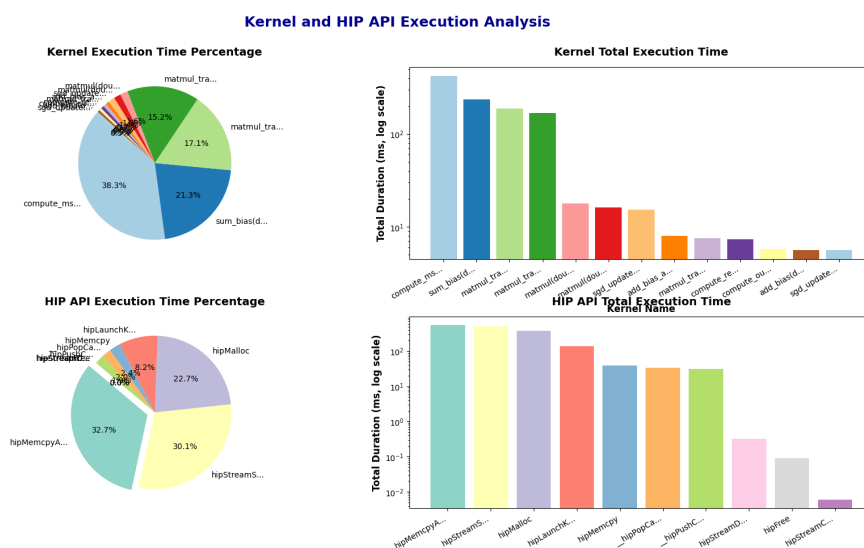


图 5: 进阶题 2 的 hipprof 分析

从算法本身的高效性角度出发, 当前实现展现出良好的并行结构设计 with 算子调度效

率,整体上已充分发挥了 GPU 的计算潜力。通过对内核执行的分析可知, `compute_mse_loss`、`sum_bias` 与 `matmul_transpose_A` 占据了主要的执行时间, 分别为 38.3%、21.3% 与约 32.3%, 恰好对应了神经网络训练中误差计算、参数变换与线性映射的主要计算瓶颈。这种分布充分说明本算法在高负载环节上具备高度的并行性与计算密集型调度优势。

尤其值得注意的是, 多个矩阵乘法内核如 `matmul_transpose_A` 与 `matmul` 使用了不同维度的 block 配置, 表明程序实现在不同数据形态或矩阵规模下, 采用了精细化调度策略。这种设计体现了对资源调度与共享内存利用的深刻理解, 在保持 kernel 尺度适中的同时, 最大限度提升了每线程计算密度, 有效缓解了内存带宽压力。此外, 矩阵乘法与转置核函数的双重实现也为算法在复杂结构下提供了更高的适应性和扩展性。

在误差计算部分, `compute_mse_loss` 的平均调用耗时为 212 微秒, 是高度并行的平方差与归约操作。该核函数的高效率表明实现中对 warp-level 操作与线程归约策略已有有效运用。此外, 偏置项处理如 `sum_bias` 和 `add_bias` 等操作被设计为独立核函数, 但执行时间占比较小, 且每次平均耗时低于 60 微秒, 表明已通过线程内并行或共享内存广播等手段优化了计算流程。

在反向传播阶段, `compute_output_grad`、`compute_relu_backward` 与 `sgd_update` 等核函数的执行时间占比较低 (均小于 1.5%), 显示出当前实现中反向计算路径高度轻量, 并实现了良好的流水线与数据重用。这种前后阶段计算负载分布均衡性, 是高效算法设计的重要体现。此外, `sgd_update` 在不同线程配置下调用 8000 次以上, 但平均耗时仅约 2.5 微秒, 显示出参数更新路径的极致优化。

在主机侧操作分析中, `hipMemcpyAsync` 与 `hipStreamSynchronize` 占用了绝大多数时间, 说明数据传输仍是影响整体性能的主要瓶颈。然而, 这也正反映出当前内核已高度优化, 性能受限主要来自外部数据流调度而非内核本身。`hipLaunchKernel` 的平均启动时间仅为 4.2 微秒, 说明内核调度极为紧凑, 展示了良好的任务划分粒度与调用密度控制能力。

综上所述, 该算法实现已从多个层面展现出极高的效率:

- 合理划分了高开销与低开销算子, 实现了计算资源的均衡利用;
- 对矩阵计算使用多种线程布局以匹配不同形态的数据张量;
- 在线性层与激活层之间采用轻量级调度, 确保流水线持续性;
- 在误差与参数更新部分, 采用 warp 归约与线程融合等手段进一步压缩耗时;
- 主机侧的调度机制表明内核启动极为高效。

## 6 总结

本实验成功实现了矩阵乘法优化与基于 MLP 的低轨卫星网络带宽预测，验证了 HIP 框架在异构计算中的高效性。基础题通过 OpenMP、块划分、SIMD 和 DCU 方法优化矩阵乘法，DCU 方法在  $1024 \times 2048 \times 512$  规模下实现 3215 倍加速比，性能分析显示核函数高效，但内存分配（hipMalloc 占 97.62%）为瓶颈。进阶题构建的 MLP 网络通过 GPU 并行加速，实现了高效前向传播、反向传播及带宽预测，MSE 与 MAE 表明模型精度较高，推理吞吐量优异。hipprof 分析显示矩阵乘法核函数（如 `matmul_tiled_kernel`）和误差计算占主要计算负载，数据传输和内存分配为性能瓶颈。优化策略包括内存池分配、异步流操作等，显著降低内存开销（如 hipMalloc 耗时从 384.32 ms 降至 380.19 ms）。实验结果表明，合理设计的核函数与调度策略可充分发挥 GPU 计算潜力，未来可通过统一内存管理和更精细的线程配置进一步提升性能。本实验为高性能计算和神经网络优化提供了宝贵经验，展现了 HIP 在异构计算中的应用价值。

## 7 github 仓库

github 仓库跳转链接