



多核程序设计

作业二：计算高维空间中的最近邻

数据科学与计算机学院 17大数据与人工智能
17341015 陈鸿峥

目录

1	题目描述	2
2	问题简答	2
3	具体实现	2
3.1	CPU版本	3
3.2	GPU基线版本	4
3.3	GPU共享内存版本	6
3.4	GPU并行归约版本	7
3.5	GPU Warp原语版本	8
3.6	GPU Block归约版本	10
4	实验设置与结果	11
4.1	环境设置	11
4.2	实验结果	11
A	完整实验数据	12

一、 题目描述

计算高维空间中的最近邻(nearest neighbor)

- 输入：查询点集，参考点集，空间维度 k
- 输出：对每个查询点，输出参考点集中最近邻的序号

回答以下问题：

1. 介绍程序整体逻辑，包含的函数，每个函数完成的内容。（10分）
 - 对于核函数，应该说明每个线程块及每个线程所分配的任务
2. 解释程序中涉及哪些类型的存储器（如全局内存、共享内存等），并通过分析数据的访存模式及该存储器的特性说明为何使用该种存储器。（15分）
3. 针对查询点集为1个点及1024个点给出两个版本，并说明设计逻辑（例如任务分配方式）的异同。如两个版本使用同一逻辑，需说明原因。（35分）
4. 请给出一个基础版本（baseline）及至少一个优化版本。并分析说明每种优化对性能的影响。（40分）
 - 具体得分根据优化难度（技巧及工作量）而定
5. 选做：使用空间划分数据结构加速查询（如KD-Tree, BVH-Tree）。（20分）

二、 问题简答

1. 参见第3节，每个线程所做的任务均用加粗字体标注。
2. 参见第3节GPU共享内存、Warp原语版本等，不同类型的存储器均使用不同颜色字体进行标注。
3. 参见第3节，除了GPU Block归约版本，其他GPU优化版本对于查询点集为1个点和1024个点的任务分配方式均相同，但具体每个线程做的任务有所不同，可见下文加粗字体。基本的分配策略都是开设 m 个Block给 m 个查询点使用，而每个Block内的每个线程则用于计算每个参考点到该查询点的距离。这样的好处在于查询点的坐标可以直接存储在**共享内存**中，供所有参考点使用；并且所有的并行归约操作都限制在一个Block内，可以直接通过**共享内存**传输数据，而无需通过全局内存，大大降低了访存的开销。
4. 参见第3节，一共实现了五个版本的GPU程序进行优化与比较，实验结果见第4节。

三、 具体实现

本次实验我总共实现了六个版本的程序，包括CPU版本、GPU基线版本、GPU共享内存版本、GPU并行归约版本、GPU Warp原语版本，GPU Block归约版本。

1. CPU版本

最基础的实现即助教提供的CPU版本程序，依次遍历所有的查询点，对于每一个查询点，遍历所有的参考点，并（遍历每个维度）计算每个查询点和参考点之间的距离。从这些距离中选取最小的，其下标即为所求的最近邻。

```

1  /*!
2  * Naive CPU implementation
3  * used to test the correctness of the results
4  * \param k The dimension size of the points
5  * \param m The nubmer of search points
6  * \param n The number of reference points
7  * \param searchPoints
8  * \param referencePoints
9  * \param results
10 * \return void. Results will be put in result.
11 */
12 extern void cudaCallbackCPU(int k, int m, int n, float *searchPoints,
13                             float *referencePoints, int **results) {
14
15     int *tmp = (int*)malloc(sizeof(int)*m);
16     int minIndex;
17     float minSquareSum, diff, squareSum;
18
19     // Iterate over all search points
20     for (int mInd = 0; mInd < m; mInd++) {
21         minSquareSum = -1;
22         // Iterate over all reference points
23         for (int nInd = 0; nInd < n; nInd++) {
24             squareSum = 0;
25             for (int kInd = 0; kInd < k; kInd++) {
26                 diff = searchPoints[k*mInd+kInd] - referencePoints[k*nInd+kInd];
27                 squareSum += (diff * diff);
28             }
29             if (minSquareSum < 0 || squareSum < minSquareSum) {
30                 minSquareSum = squareSum;
31                 minIndex = nInd;
32             }
33         }
34         tmp[mInd] = minIndex;
35     }
36
37     *results = tmp;
38     // Note that you don't have to free searchPoints, referencePoints, and
39     // *results by yourself

```

40 }

2. GPU基线版本

GPU基线(baseline)版本即将CPU版本的最外层循环进行并行。每个线程处理一个查询点，同样是遍历所有的参考点并计算查询点与参考点之间的距离。每个线程都各自计算出所有参考点的最小距离，并将结果直接写入全局内存中。核函数如下所示。

```

1  /*!
2  * Core execution part of CUDA
3  * that calculates the nearest neighbor of each search point.
4  * \param k The dimension size of the points
5  * \param m The nubmer of search points
6  * \param n The number of reference points
7  * \param searchPoints
8  * \param referencePoints
9  * \param output
10 * \return void. Results will be put in output
11 */
12 __global__ void kernel(int k, int m, int n, float* searchPoints, float*
    ↪ referencePoints, int* output) {
13     int pid = threadIdx.x;
14     int minIdx;
15     float minSquareSum = -1;
16     float diff, squareSum;
17     // Iterate over all reference points
18     if (pid < m) {
19         for (int nInd = 0; nInd < n; nInd++) { // ref points
20             squareSum = 0;
21             for (int kInd = 0; kInd < k; kInd++) { // dimension
22                 diff = searchPoints[k * pid + kInd]
23                     - referencePoints[k * nInd + kInd];
24                 squareSum += (diff * diff);
25             }
26             if (minSquareSum < 0 || squareSum < minSquareSum) {
27                 minSquareSum = squareSum;
28                 minIdx = nInd;
29             }
30         }
31         output[pid] = minIdx;
32     }
33 }
```

GPU调用的接口如下所示，主要分为几个步骤：

1. 为查询点坐标、参考点坐标、输出位置分配GPU内存

2. 将查询点坐标、参考点坐标从主机端(CPU)传到设备端(GPU)
3. 调用核函数, 可以看到开设了1个Block, 并且有 m 个线程。
4. 将结果从GPU运回CPU
5. 释放GPU内存

```

1  /*!
2   * Wrapper of the CUDA kernel
3   * used to be called in the main function
4   * \param k The dimension size of the points
5   * \param m The nubmer of search points
6   * \param n The number of reference points
7   * \param searchPoints
8   * \param referencePoints
9   * \param results
10  * \return void. Results will be put in result.
11  */
12  extern void cudaCallbackGPU_baseline(int k, int m, int n, float *searchPoints,
13                                     float *referencePoints, int **results) {
14      float *searchPoints_d, *referencePoints_d;
15      int* output_d;
16
17      // Allocate device memory and copy data from host to device
18      CHECK(cudaMalloc((void **)&searchPoints_d, sizeof(float)*m*k));
19      CHECK(cudaMalloc((void **)&referencePoints_d, sizeof(float)*n*k));
20      CHECK(cudaMalloc((void **)&output_d, sizeof(int)*m));
21      CHECK(cudaMemcpy(searchPoints_d, searchPoints, sizeof(float)*m*k,
22                      ↪ cudaMemcpyHostToDevice));
23
24      CHECK(cudaMemcpy(referencePoints_d, referencePoints, sizeof(float)*n*k,
25                      ↪ cudaMemcpyHostToDevice));
26
27      // Invoke the device function
28      kernel<<< 1, m >>>(k, m, n, searchPoints_d, referencePoints_d, output_d);
29      cudaDeviceSynchronize();
30
31      // Copy back the results and de-allocate the device memory
32      *results = (int *)malloc(sizeof(int)*m);
33      assert(results != NULL);
34      CHECK(cudaMemcpy(*results, output_d, sizeof(int)*m, cudaMemcpyDeviceToHost));
35
36      // int *cpu_results;
37      // cudaCallbackCPU(k, m, n, searchPoints, referencePoints, &cpu_results);
38      // for (int i = 0; i < m; ++i)
39      //     assert(cpu_results[i] == (*results)[i]);

```

```

38     CHECK(cudaFree(searchPoints_d));
39     CHECK(cudaFree(referencePoints_d));
40     CHECK(cudaFree(output_d));
41 }

```

注释中的部分最开始用于核验CPU版本的程序与GPU版本程序的正确性。

3. GPU共享内存版本

注意到由于每次计算距离时都会涉及到查询点集的大量内存访问，因此一个优化思路即将查询点集放置在共享内存中；至于参考点集则仍放置在全局内存中，因为每个点坐标只会被访问一次，将其读取至共享内存也不会减少访问时间。

因此这里每个Block计算一个查询点的最近邻，每个Thread分别计算 q 个参考点与查询点的距离。由于使用的GPU每个Block的最大线程数为1024，故 $q = \lceil m/1024 \rceil$ 。具体流程如下：

- 每个Block先开取 k 个线程，读取查询点的 k 维坐标向量到共享内存s_mem中
- 创建共享内存dist和dist_idx分别存储该Block内的最小距离值及其下标。注意到每个线程处理 q 个参考点，因此dist和dist_idx中的元素对应的也是这 q 个点的最小距离值及其下标（由于同个Block中的所有线程处理的是同个查询点，因此距离数组无需在Block之间传递，而可以采用共享内存方式实现）。
- 使用单一线程进行归约操作，对dist数组进行遍历，寻找其中的最小值
- 每个Block将最终结果（距离最小值对应坐标）写入全局数组output中

完整核函数代码如下所示。

```

1  __global__ void kernel_sharedmem(int k, int m, int n, float* searchPoints, float*
    ↪ referencePoints, int* output) {
2      int bid = blockIdx.x; // each block, one search point
3      int tid = threadIdx.x;
4      int n_points = ((n % 1024 == 0) ? n / 1024 : n / 1024 + 1);
5      int searchId = bid;
6      int referenceId = tid; // [tid,tid+n_points]
7      int minIdx;
8      float diff, squareSum;
9      __shared__ float s_mem[16];
10     if (tid < k) {
11         s_mem[tid] = searchPoints[k * searchId + referenceId];
12     }
13     __syncthreads();
14     __shared__ float dist[1024];
15     __shared__ int dist_idx[1024];
16     float minSquareSum = -1;

```

```

17     for (int i = 0; i < n_points; ++i) {
18         squareSum = 0;
19         int refId = referenceId * n_points + i;
20         for (int kInd = 0; kInd < k; kInd++) { // dimension
21             diff = s_mem[kInd] - referencePoints[k * refId + kInd];
22             squareSum += (diff * diff);
23         }
24         if (minSquareSum < 0 || squareSum < minSquareSum) {
25             minSquareSum = squareSum;
26             minIdx = refId;
27         }
28     }
29     dist[referenceId] = minSquareSum;
30     dist_idx[referenceId] = minIdx;
31     __syncthreads();
32     if (referenceId == 0) {
33         float minSquareSum = -1;
34         int bound = min(n, 1024);
35         for (int i = 0; i < bound; ++i) {
36             squareSum = dist[i];
37             if (minSquareSum < 0 || squareSum < minSquareSum) {
38                 minSquareSum = squareSum;
39                 minIdx = dist_idx[i];
40             }
41         }
42         output[searchId] = minIdx;
43     }
44 }

```

4. GPU并行归约版本

在上述**共享内存**版本的基础上，添加并行归约的功能。需要在1024个元素中找最小值，这个可以采用树形的方式进行归约，如图1所示。

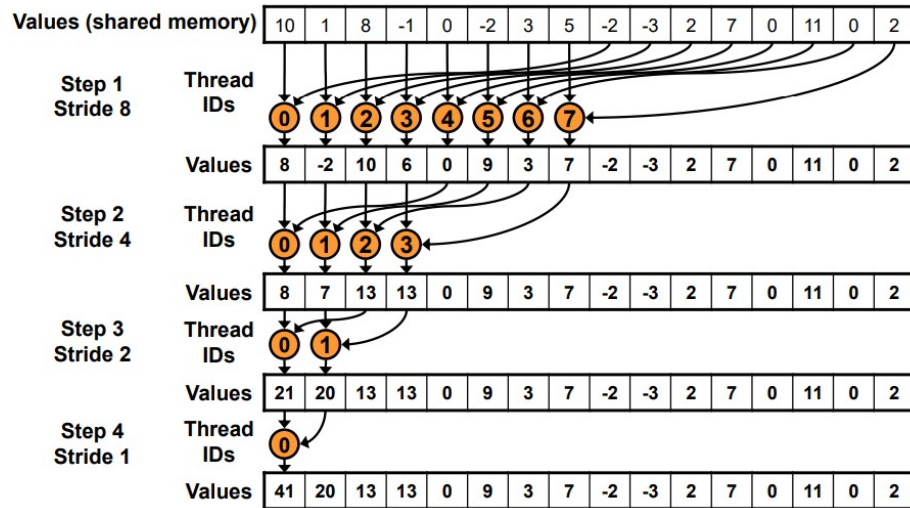


图 1: 树形方式归约

注意这里读取的都是连续内存，通过加步长来消除存储体冲突。每个线程计算两个元素间的最小值，最终1024个元素的最小值会被存放在数组第一个元素中。由于全程的操作都在**共享内存**中完成，因此速度非常快。最后只有一次全局内存的写操作，即将最终结果写入output中。

核函数并行归约部分代码如下所示。

```

1  // parallel reduction
2  // only support power of 2 at this time
3  for (int stride = 512; stride > 0; stride >>= 1) {
4      if (n > stride && tid < stride) {
5          if (dist[tid] > dist[tid + stride]) {
6              dist[tid] = dist[tid + stride];
7              dist_idx[tid] = dist_idx[tid + stride];
8          }
9      }
10     __syncthreads();
11 }
12 if (tid == 0)
13     output[searchId] = dist_idx[0];

```

为确保参考点数目小于1024时也能够正常执行，这里的判断条件还添加了对 n 与步长关系的判断。不过这个程序一个小缺陷在于只能支持2的指数次幂的归约。

5. GPU Warp原语版本

这里同样是对归约部分进行优化，由于采用树形归约方式，当线程数少于32时，就只剩第一个线程束(warp)还在工作了。因此这时候可以直接采用线程束内的并行机制，在课件中可以通过展开最后一个线程束而无需添加__syncthreads来实现并行归约求和。但是对于求最小值

这种涉及到判断语句的归约，则没有办法简单通过上述方法实现自动同步，否则会导致结果运算错误。

故解决方案是采用warp层面的原语操作，如图2所示，CUDA 9引入了几个新的原语，使得并行归约操作变得更加高效便捷。对于一个线程束内的32个线程，`__shfl_down_sync`可以实现线程之间的数据传输直接通过寄存器进行。通过指定掩码、传输变量及偏移量，即可在线程束内实现快速高效的数据交换，进而达成归约的目的。

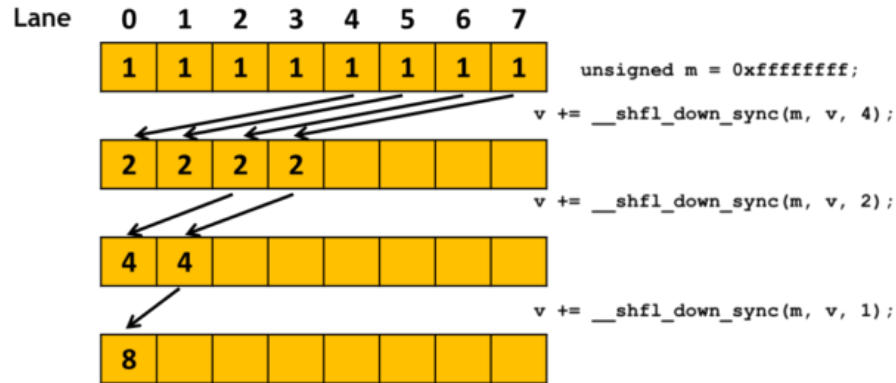


图 2: `__shfl_down_sync`原语操作

具体到本问题来说，我将最后一个线程束进行包裹，实现了`warpReduceMin`函数，同样是树形归约，但直接将高位线程数据搬移至低位线程求最小值，而不用再计算对应数组地址。

```

1  __inline__ __device__ void warpReduceMin(float& val, int& idx)
2  {
3      for (int stride = warpSize / 2; stride > 0; stride >>= 1) {
4          float tmpVal = __shfl_down_sync(FULL_MASK, val, stride);
5          int tmpIdx = __shfl_down_sync(FULL_MASK, idx, stride);
6          if (tmpVal < val) {
7              val = tmpVal;
8              idx = tmpIdx;
9          }
10     }
11 }
```

注意到这里无需显式同步，而且数据直接通过寄存器交换，速度也会大幅提升。

核函数部分修改为下面代码。

```

1  // parallel reduction
2  // only support power of 2 at this time
3  for (int stride = 512; stride >= 32; stride >>= 1) {
4      if (n > stride && tid < stride) {
5          if (dist[tid] > dist[tid + stride]) {
6              dist[tid] = dist[tid + stride];

```

```

7         dist_idx[tid] = dist_idx[tid + stride];
8     }
9 }
10 __syncthreads();
11 }
12 if (tid < 32)
13     warpReduceMin(dist[tid], dist_idx[tid]);
14 if (tid == 0)
15     output[searchId] = dist_idx[0];

```

6. GPU Block归约版本

有了warp内的归约求最值，更进一步，可以将一个block内的所有warp都按照这种方式实现，这样将大幅减少共享内存的访问，同时也减少线程之间的同步。

故有下列blockReduceMin函数，如果一个block内有1024个线程，那么将有32个warp。对每个warp先调用warpReduceMin求出每个warp的最小值，然后再对这32个最小值求全局最小。最后0号线程的结果即为全局最小值。

```

1 __inline__ __device__ void blockReduceMin(float& val, int& idx)
2 {
3     static __shared__ float values[32];
4     static __shared__ int indices[32];
5     int lane = threadIdx.x % warpSize;
6     int wid = threadIdx.x / warpSize;
7
8     warpReduceMin(val, idx); // partial reduction for each warp in a block
9
10    if (lane == 0) { // write back to shared mem
11        values[wid] = val;
12        indices[wid] = idx;
13    }
14
15    __syncthreads();
16
17    // read from shared memory only if that warp existed
18    if (threadIdx.x < blockDim.x / warpSize) {
19        val = values[lane];
20        idx = indices[lane];
21    } else {
22        val = INT_MAX;
23        idx = 0;
24    }
25
26    if (wid == 0) {
27        warpReduceMin(val, idx); // final reduce within first warp

```

```
28     }  
29 }
```

注意如果参考点集小于1024个，则会回滚回Warp原语的版本，即只对最后一个线程束进行归约求最小值。

四、实验设置与结果

1. 环境设置

为了避免占用学院集群环境的大量资源，我采用了另外一台服务器进行测试。该服务器上装备了2个Intel Xeon Gold 5118处理器（24物理核/48逻辑核），另有一块Titan V GPU。操作系统为Ubuntu 18.04 LTS，编译器采用gcc 7.4.0和nvcc 10.2。同时，我也在学院集群上进行了简单测试，运行时间均与下述在服务器上的测试差异不大。

所有版本的正确性均已测试通过，下列所有实验均为运行多次取平均的结果。

2. 实验结果

实验结果如图3所示，可以看到CPU版本的程序在样例较小时表现都比较好，但是一旦查询/参考点的数目提升了，其性能就下降得很快，与GPU优化后的版本差了2个数量级。

另外，可以看到**共享内存**的使用给GPU的性能带来了极大的提升，至少是1个数量级的优化。GPU Baseline第一组数据的突起是由于冷启动造成的，后面GPU程序的运行则没有这种问题。

由于本次实验提供的测试样例规模依然较小，因此采用了不同归约优化方式，性能提升在图中看并不明显。但是从完整实验数据（见附录A）中可以看出无论是并行归约，还是采用warp原语，对GPU性能的影响依然是非常正面的。

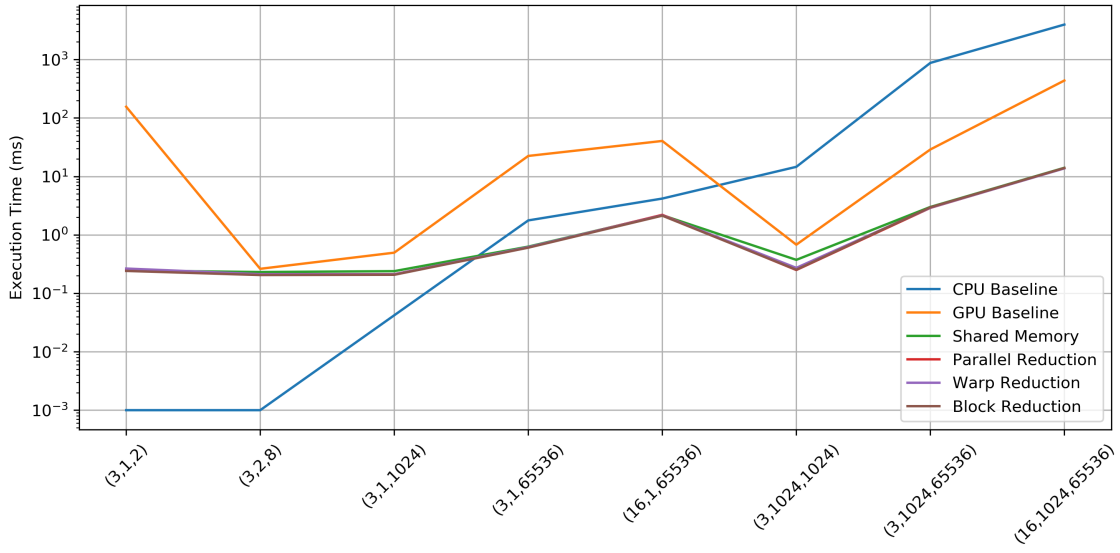


图 3: CUDA不同优化版本比较

附录 A. 完整实验数据

表 1: 不同规模数据下不同版本程序的运行时间(ms)

	(3,1,2)	(3,2,8)	(3,1,1024)	(3,1,65536)	(16,1,65536)	(3,1024,1024)	(3,1024,65536)	(16,1024,65536)
CPU基线	0.001	0.001	0.042	1.764	4.171	14.559	872.909	3955.051
GPU基线	155.364	0.262	0.495	22.362	40.349	0.678	28.838	436.137
GPU共享内存	0.245	0.231	0.239	0.628	2.168	0.373	2.997	14.024
GPU并行归约	0.253	0.212	0.212	0.608	2.196	0.250	2.891	13.815
GPU Warp原语	0.265	0.212	0.211	0.616	2.149	0.272	2.909	13.696
GPU Block归约	0.241	0.205	0.207	0.603	2.126	0.253	2.956	13.892