

HW3

112062706 林泳成

- HW3
 - Implementation
 - HW3-1
 - HW3-2、HW3-3
 - Profiling Results
 - Experiment & Analysis
 - Optimization
 - Time Distribution
 - Conclusion

Implementation

HW3-1

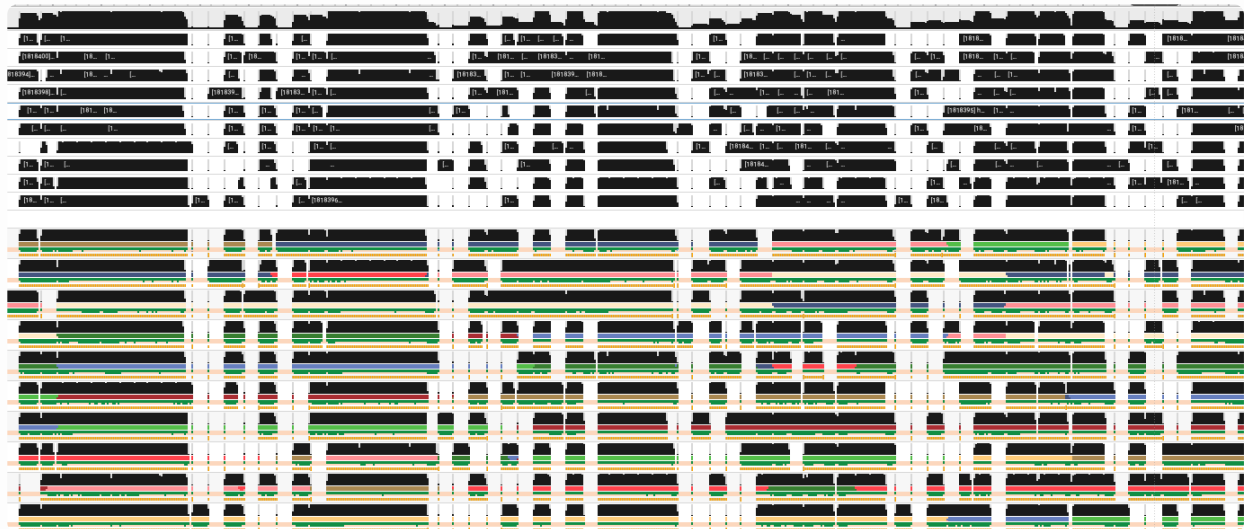
```

void cal(
    int B, int Round, int block_start_x, int block_start_y, int block_width, int block_height) {
#pragma omp parallel num_threads(NUM_THREADS)
{
    int block_end_x = block_start_x + block_height;
    int block_end_y = block_start_y + block_width;
#pragma omp for schedule(dynamic) collapse(2)
    for (int b_i = block_start_x; b_i < block_end_x; ++b_i) {
        for (int b_j = block_start_y; b_j < block_end_y; ++b_j) {
            // To calculate B*B elements in the block (b_i, b_j)
            // For each block, it need to compute B times
            int block_internal_start_x = b_i * B;
            int block_internal_end_x = (b_i + 1) * B > n ? n : (b_i + 1) * B;
            int block_internal_start_y = b_j * B;
            int block_internal_end_y = (b_j + 1) * B > n ? n : (b_j + 1) * B;

            // if (block_internal_end_x > n) block_internal_end_x = n;
            // if (block_internal_end_y > n) block_internal_end_y = n;
            for (int k = Round * B; k < (Round + 1) * B && k < n; ++k) {
                // To calculate original index of elements in the block (b_i, b_j)
                // For instance, original index of (0,0) in block (1,2) is (2,5) for V=6,B=2
                for (int i = block_internal_start_x; i < block_internal_end_x; ++i) {
                    int dik = Dist[i][k];
                    if (dik == INF) continue;
#pragma omp reduction(+:Dist)
                    for (int j = block_internal_start_y; j < block_internal_end_y; ++j) {
                        if (Dist[k][j] != INF && dik + Dist[k][j] < Dist[i][j]) {
                            Dist[i][j] = dik + Dist[k][j];
                        }
                    }
                }
            }
        }
    }
}
}

```

原本是用openMP直接對blocked floyd warshall做修改。但是好像是每個子工作太大，沒辦法做到太fine grained



用nsys看執行狀況發現等待的情況很嚴重，應該是子工作太大導致的
 這個方法因為SIMD包住的範圍包含了三個額外的for迴圈做陣列更新，即便是用
 schedule(dynamic)，也會因為一個thread跑一次子工作太久導致所有thread都在等待
 每次結束SIMD的迴圈時，都會額外有一個process獨自執行，其他process都沒工作。好像是因為
 openMP在做完SIMD操作後會附上一個barrier，之後會由主程式做資源管理，負責在runtime做資源

回收和調度

所以後面優化就是想優化thread資源分配的overhead

如果只在程式一開始初始化一次thread就可以避免中間執行花太多時間在資源管理

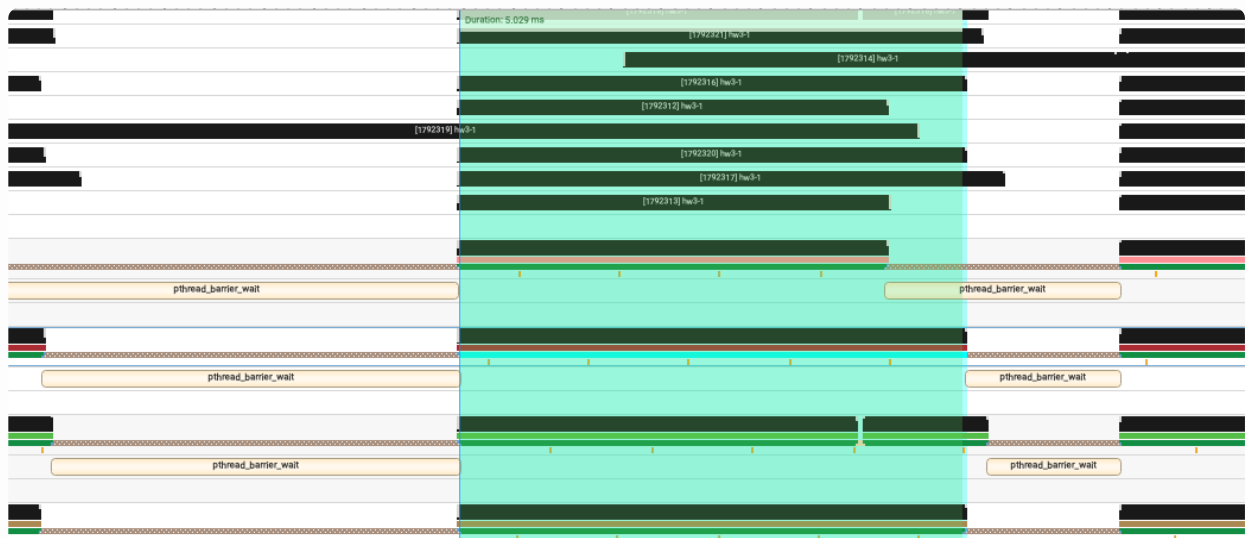
優化後

```
int jobPerThread = ceil(n * n, NUM_THREADS);

curX = *tid * jobPerThread;
for(int i = 0; i < jobPerThread; ++i){
    if(Dist[curX/n][curX%n] > Dist[curX/n][k] + Dist[k][curX%n])
        Dist[curX/n][curX%n] = Dist[curX/n][k] + Dist[k][curX%n];
    ++curX;
}
pthread_barrier_wait(&barrier);
```

由於之前的做法工作沒辦法切得很好，會有空隙，所以想說用pthread在每一輪就先把工作切得剛剛好，每個thread在每一輪就是處理 $\text{ceil}(n*n, \text{NUM_THREAD})$ 的工作量

結果表現好像沒有太好，其中一個部分是每個thread被fork出來的時間是由for loop一個一個分配出來的，所以先執行的thread就需要等待較慢初始化的thread



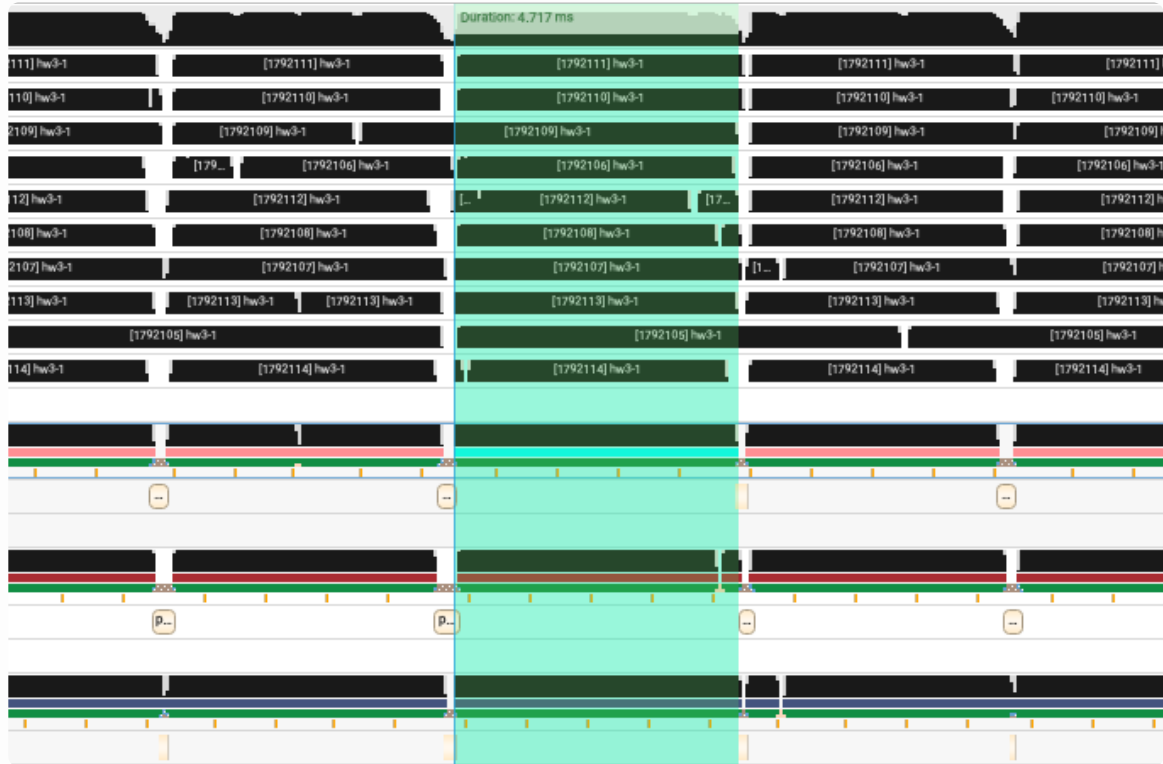
每一輪計算都會花5毫秒

```
for(k = 0; k < n; ++k){
    for(i = *tid; i < n; i += NUM_THREADS){
        tmpDist = Dist[i][k];
        for(j = 0; j < n; ++j){
            if(Dist[i][j] > tmpDist + Dist[k][j]){
                Dist[i][j] = tmpDist + Dist[k][j];
            }
        }
    }
    pthread_barrier_wait(&barrier);
}
```

因為之前的做法每一輪的計算會花很多時間，所以後來去看了cache要怎麼存取才會快

就發現主要是cache miss花太多時間，就連之後又測試了__builtin_prefetch都沒辦法在不同thread之間maintain cache，導致整體access時間增加

所以就改成每個thread同時只能抓一行的陣列來做運算，結果大幅改善cache miss的問題



每一輪計算只花4.7毫秒

HW3-2、HW3-3

```
for (int r = 0; r < round; ++r) {
    printf("%d %d\n", r, round);
    fflush(stdout);
    /* Phase 1*/
    cal1<<<gridsPerBlock1, threadsPerBlock>>>(r, deviceDist, pitch / sizeof(int));

    /* Phase 2*/
    cal2<<<gridsPerBlock2, threadsPerBlock>>>(r, deviceDist, pitch / sizeof(int));

    /* Phase 3*/
    cal3<<<gridsPerBlock3, threadsPerBlock>>>(r, deviceDist, pitch / sizeof(int));
}
```

```
for (int r = 0; r < round; ++r) {
    // 擁有這一輪需要的row的device釋出這個row給另一個device
    if (r >= yStart && r < yStart + gridsPerBlock3.y) {
        cudaMemcpy2DAsync(hostDist + r * B * pad_n, pad_n * sizeof(int), (char*)deviceDist[omp_get_thread_num()] + r * B * pitch, pitch, pad_n * sizeof(int), B, cudaMemcpyDeviceToHost);
    }
    cudaStreamSynchronize(stream);
    #pragma omp barrier

    if (r < yStart || r >= yStart + gridsPerBlock3.y) {
        cudaMemcpy2DAsync((char*)deviceDist[omp_get_thread_num()] + r * B * pitch, pitch, hostDist + r * B * pad_n, pad_n * sizeof(int), pad_n * sizeof(int), B, cudaMemcpyHostToDevice);
    }
    cal1<<<gridsPerBlock1, threadsPerBlock>>>(r, deviceDist[omp_get_thread_num()], pitch / sizeof(int));
    cal2<<<gridsPerBlock2, threadsPerBlock>>>(r, deviceDist[omp_get_thread_num()], pitch / sizeof(int));
    cal3<<<gridsPerBlock3, threadsPerBlock>>>(r, deviceDist[omp_get_thread_num()], pitch / sizeof(int), yStart);
}
```

資料的分配上都是照著每個phase要計算的工作去切

只有在HW3-3是把整個母工作分成上下兩大塊個別計算

HW3-2和HW3-3我都是選擇blocking factor 64、block裡的thread都是(32X32)。

Blocking factor選64是因為一個block是32X32個thread，而且share memory又可以長寬各開到64X64，所以在嘗試32和64兩種blocking factor後選擇效能比較好的64。

至於block的thread分配則是選擇講義裡建議的x, y各是32倍數。

至於grid的分布則是把graph的長寬除以block大小。

```
Device 0: "NVIDIA GeForce GTX 1080"
CUDA Driver Version / Runtime Version      12.6 / 12.6
CUDA Capability Major/Minor version number: 6.1
Total amount of global memory:              8107 MBytes (8500871168 bytes)
(20) Multiprocessors, (128) CUDA Cores/MP:  2560 CUDA Cores
GPU Max Clock rate:                         1835 MHz (1.84 GHz)
Memory Clock rate:                          5005 Mhz
Memory Bus Width:                           256-bit
L2 Cache Size:                              2097152 bytes
Maximum Texture Dimension Size (x,y,z)      1D=(131072), 2D=(131072, 65536), 3D=(16384, 16384, 16384)
Maximum Layered 1D Texture Size, (num) layers 1D=(32768), 2048 layers
Maximum Layered 2D Texture Size, (num) layers 2D=(32768, 32768), 2048 layers
Total amount of constant memory:             65536 bytes
Total amount of shared memory per block:     49152 bytes
Total number of registers available per block: 65536
Warp size:                                   32
Maximum number of threads per multiprocessor: 2048
Maximum number of threads per block:         1024
Max dimension size of a thread block (x,y,z): (1024, 1024, 64)
Max dimension size of a grid size    (x,y,z): (2147483647, 65535, 65535)
Maximum memory pitch:                       2147483647 bytes
Texture alignment:                           512 bytes
Concurrent copy and kernel execution:        Yes with 2 copy engine(s)
Run time limit on kernels:                   No
Integrated GPU sharing Host Memory:          No
Support host page-locked memory mapping:     Yes
Alignment requirement for Surfaces:          Yes
Device has ECC support:                      Disabled
Device supports Unified Addressing (UVA):     Yes
Supports Cooperative Kernel Launch:          Yes
Supports MultiDevice Co-op Kernel Launch:    Yes
```

```
#pragma omp parallel num_threads(2)
{
    cudaSetDevice(omp_get_thread_num());
    cudaStream_t stream;
    cudaStreamCreate(&stream);
    dim3 gridsPerBlock3(round, (omp_get_thread_num() == 1) && (blockNum & 1) ? ceil(blockNum, 2) : blockNum / 2);
    // y座標從第幾個block開始
    int yStart = blockNum / 2 * omp_get_thread_num();
    cudaMallocPitch(&deviceDist[omp_get_thread_num()], &pitch, pad_n * sizeof(int), pad_n);
    // cudaMemcpy2DAsync(deviceDist[omp_get_thread_num()], pitch, hostDist, pad_n * sizeof(int), pad_n * sizeof(int), pad_n, cudaMemcpyHostToDevice, stream);
    cudaMemcpy2DAsync((char*)deviceDist[omp_get_thread_num()] + yStart * B * pitch, pitch, hostDist + yStart * B * pad_n, pad_n * sizeof(int), pad_n * sizeof(int), gridsPerBlock3.y * B, cudaMemcpyHostToDevice, stream);
    cudaStreamSynchronize(stream);
    for (int r = 0; r < round; ++r) {
        // 擁有這一輪需要的row的device釋出這個row給另一個device
        if (r >= yStart && r < yStart + gridsPerBlock3.y) {
            cudaMemcpy2DAsync(hostDist + r * B * pad_n, pad_n * sizeof(int), (char*)deviceDist[omp_get_thread_num()] + r * B * pitch, pitch, pad_n * sizeof(int), B, cudaMemcpyDeviceToHost, stream);
        }
        cudaStreamSynchronize(stream);
        #pragma omp barrier

        if (r < yStart || r >= yStart + gridsPerBlock3.y) {
            cudaMemcpy2DAsync((char*)deviceDist[omp_get_thread_num()] + r * B * pitch, pitch, hostDist + r * B * pad_n, pad_n * sizeof(int), pad_n * sizeof(int), B, cudaMemcpyHostToDevice, stream);
        }
        cal1<<<gridsPerBlock1, threadsPerBlock>>>(r, deviceDist[omp_get_thread_num()], pitch / sizeof(int));
        cal2<<<gridsPerBlock2, threadsPerBlock>>>(r, deviceDist[omp_get_thread_num()], pitch / sizeof(int));
        cal3<<<gridsPerBlock3, threadsPerBlock>>>(r, deviceDist[omp_get_thread_num()], pitch / sizeof(int), yStart);
    }
    cudaMemcpy2DAsync(hostDist + yStart * B * pad_n, pad_n * sizeof(int), (char*)deviceDist[omp_get_thread_num()] + yStart * B * pitch, pitch, pad_n * sizeof(int), gridsPerBlock3.y * B, cudaMemcpyDeviceToHost, stream);
    cudaStreamSynchronize(stream);
    cudaStreamDestroy(stream);
}
```

HW3-3因為會把整個graph切成上下兩大塊，所以在溝通上會是根據每一輪，擁有那一round資料的人把資料丟出來給沒有的GPU拿。這樣就可以讓兩個GPU平行計算phase 3，但是phase 1因為只需要一個GPU算，phase 2因為工作量少所以效果都有限。最後再把兩個GPU上的graph組合起來就好了。

實作上有試著改變access pattern來解決bank conflict(因為phase 2會剛好用完share memory，就不能用padding)，但是效果很差。

其餘部分就是分工把整塊block搬到share memory上。

```

sharedCenter[ty * B + tx] = deviceDist[centerY * pitch + centerX];
sharedCenter[(ty+halfB) * B + tx] = deviceDist[(centerY+halfB) * pitch + centerX];
sharedCenter[ty * B + (tx+halfB)] = deviceDist[centerY * pitch + (centerX+halfB)];
sharedCenter[(ty+halfB) * B + (tx+halfB)] = deviceDist[(centerY+halfB) * pitch + (centerX+halfB)];
sharedCol[ty * B + tx] = deviceDist[colY * pitch + colX];
sharedCol[(ty+halfB) * B + tx] = deviceDist[(colY+halfB) * pitch + colX];
sharedCol[ty * B + (tx+halfB)] = deviceDist[colY * pitch + (colX+halfB)];
sharedCol[(ty+halfB) * B + (tx+halfB)] = deviceDist[(colY+halfB) * pitch + (colX+halfB)];
sharedRow[ty * B + tx] = deviceDist[rowY * pitch + rowX];
sharedRow[(ty+halfB) * B + tx] = deviceDist[(rowY+halfB) * pitch + rowX];
sharedRow[ty * B + (tx+halfB)] = deviceDist[rowY * pitch + (rowX+halfB)];
sharedRow[(ty+halfB) * B + (tx+halfB)] = deviceDist[(rowY+halfB) * pitch + (rowX+halfB)];
__syncthreads();

#pragma unroll 64
for (int k = 0; k < B; ++k) {
    // sharedCol[ty * B + tx] = min(sharedCol[ty * B + tx], sharedCol[ty * B + k] + sharedCenter[k * B + tx]);
    // sharedRow[ty * B + tx] = min(sharedRow[ty * B + tx], sharedCenter[ty * B + k] + sharedRow[k * B + tx]);
    // Update column
    if (sharedCol[ty * B + k] + sharedCenter[k * B + tx] < sharedCol[ty * B + tx]) {
        sharedCol[ty * B + tx] = sharedCol[ty * B + k] + sharedCenter[k * B + tx];
    }
    if (sharedCol[(ty+halfB) * B + k] + sharedCenter[k * B + tx] < sharedCol[(ty+halfB) * B + tx]) {
        sharedCol[(ty+halfB) * B + tx] = sharedCol[(ty+halfB) * B + k] + sharedCenter[k * B + tx];
    }
    if (sharedCol[ty * B + k] + sharedCenter[k * B + (tx+halfB)] < sharedCol[ty * B + (tx+halfB)]) {
        sharedCol[ty * B + (tx+halfB)] = sharedCol[ty * B + k] + sharedCenter[k * B + (tx+halfB)];
    }
    if (sharedCol[(ty+halfB) * B + k] + sharedCenter[k * B + (tx+halfB)] < sharedCol[(ty+halfB) * B + (tx+halfB)]) {
        sharedCol[(ty+halfB) * B + (tx+halfB)] = sharedCol[(ty+halfB) * B + k] + sharedCenter[k * B + (tx+halfB)];
    }
    // Update row
    if (sharedCenter[ty * B + k] + sharedRow[k * B + tx] < sharedRow[ty * B + tx]) {
        sharedRow[ty * B + tx] = sharedCenter[ty * B + k] + sharedRow[k * B + tx];
    }
    if (sharedCenter[(ty+halfB) * B + k] + sharedRow[k * B + tx] < sharedRow[(ty+halfB) * B + tx]) {
        sharedRow[(ty+halfB) * B + tx] = sharedCenter[(ty+halfB) * B + k] + sharedRow[k * B + tx];
    }
    if (sharedCenter[ty * B + k] + sharedRow[k * B + (tx+halfB)] < sharedRow[ty * B + (tx+halfB)]) {
        sharedRow[ty * B + (tx+halfB)] = sharedCenter[ty * B + k] + sharedRow[k * B + (tx+halfB)];
    }
    if (sharedCenter[(ty+halfB) * B + k] + sharedRow[k * B + (tx+halfB)] < sharedRow[(ty+halfB) * B + (tx+halfB)]) {

```

也有使用stream把傳資料和計算分成兩個stream，但是效果也有限。因為總共只需要傳來回一次，就要為了這來回一次創建一個stream，導致變更慢。

```

cudaStream_t stream;
cudaStreamCreate(&stream);

cudaHostRegister(hostDist, pad_n * pad_n * sizeof(int), cudaHostAllocDefault);
cudaMallocPitch(&deviceDist, &pitch, pad_n * sizeof(int), pad_n);
cudaMemcpy2DAsync(deviceDist, pitch, hostDist, pad_n * sizeof(int), pad_n * sizeof(int), pad_n, cudaMemcpyHostToDevice, stream);
// cudaMalloc(&deviceDist, pad_n * pad_n * sizeof(int));
// cudaMemcpy(deviceDist, hostDist, pad_n * pad_n * sizeof(int), cudaMemcpyHostToDevice);

cudaGetDeviceProperties(&prop, DEV_N0);
printf("maxThreadsPerBlock = %d, sharedMemPerBlock = %d\n", prop.maxThreadsPerBlock, prop.sharedMemPerBlock);

block_FW(deviceDist);
cudaMemcpy2DAsync(hostDist, pad_n * sizeof(int), deviceDist, pitch, pad_n * sizeof(int), pad_n, cudaMemcpyDeviceToHost, stream);
cudaStreamSynchronize(stream);
cudaStreamDestroy(stream);
// cudaMemcpy(hostDist, deviceDist, pad_n * pad_n * sizeof(int), cudaMemcpyDeviceToHost);

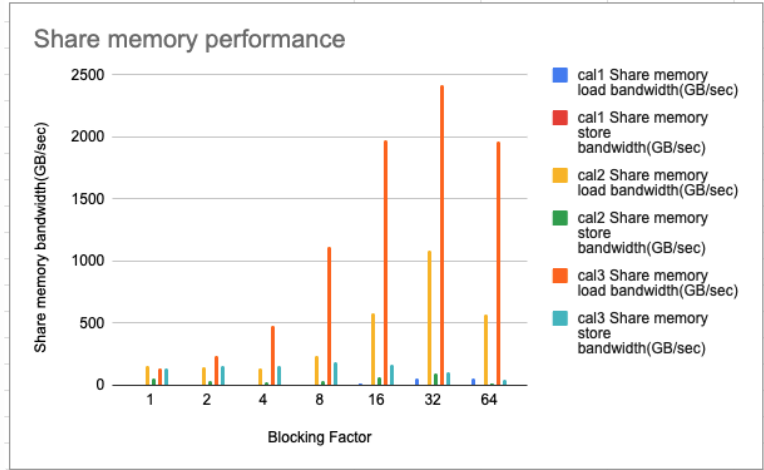
```

但是對於HW3-3多GPU來說因為每一輪都會溝通一次來回，所以多創一個stream就有價值，因此有顯著加速。

Profiling Results

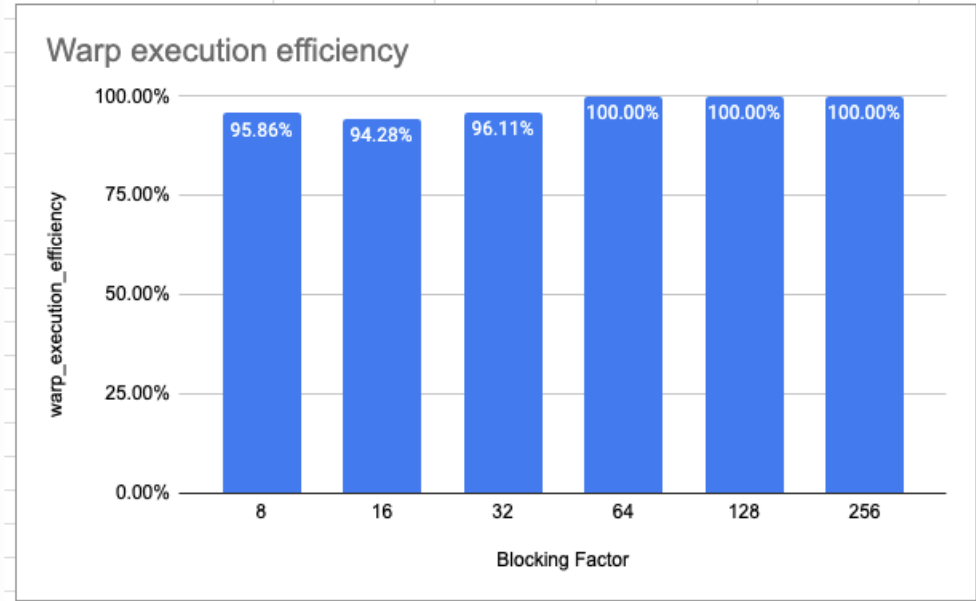
使用c10.1測資當作input，實驗不同blocking factor對GPU效能的各項指標有什麼影響

Blocking Factor	1	2	4	8	16	32	64
cal1 Share memory load bandwidth(GB/sec)	0.22088	0.37278	0.65714	2.2523	11.658	51.964	58.861
cal1 Share memory store bandwidth(GB/sec)	0.055219	0.053253	0.069555	0.24601	1.1143	4.339	1.1139
cal2 Share memory load bandwidth(GB/sec)	154.39	144.73	133.32	233.25	576.75	1086.6	567.96
cal2 Share memory store bandwidth(GB/sec)	57.897	32.563	21.383	32.109	61.984	91.783	10.355
cal3 Share memory load bandwidth(GB/sec)	140.89	241.2	475.84	1119.5	1971	2416.6	1966.7
cal3 Share memory store bandwidth(GB/sec)	140.89	160.8	158.61	186.59	164.25	100.69	40.973

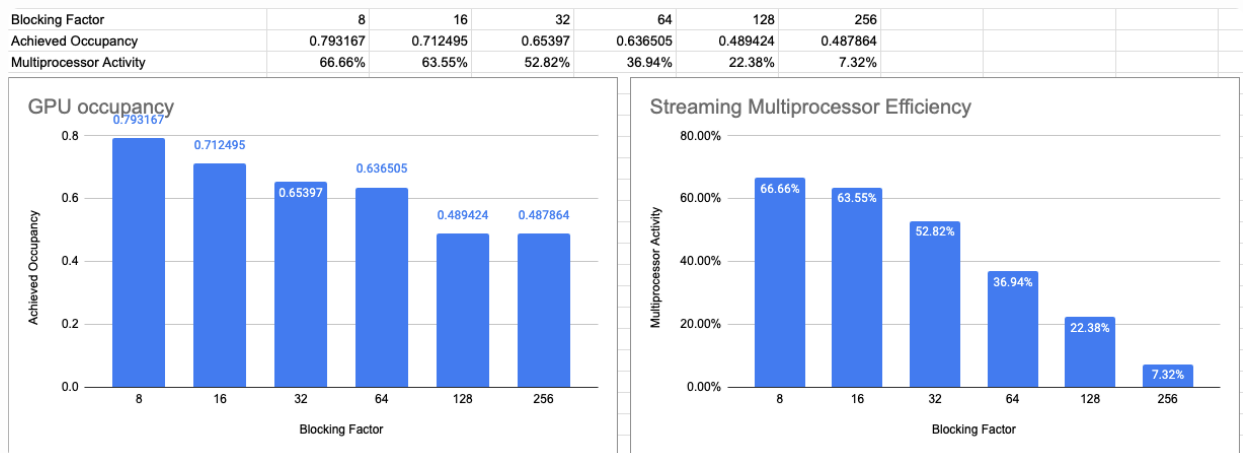


▲不同blocking factor影響GPU share memory在不同kernel call load和store操作的bandwidth

Blocking Factor	8	16	32	64	128	256
warp_execution_efficiency	95.86%	94.28%	96.11%	100.00%	100.00%	100.00%

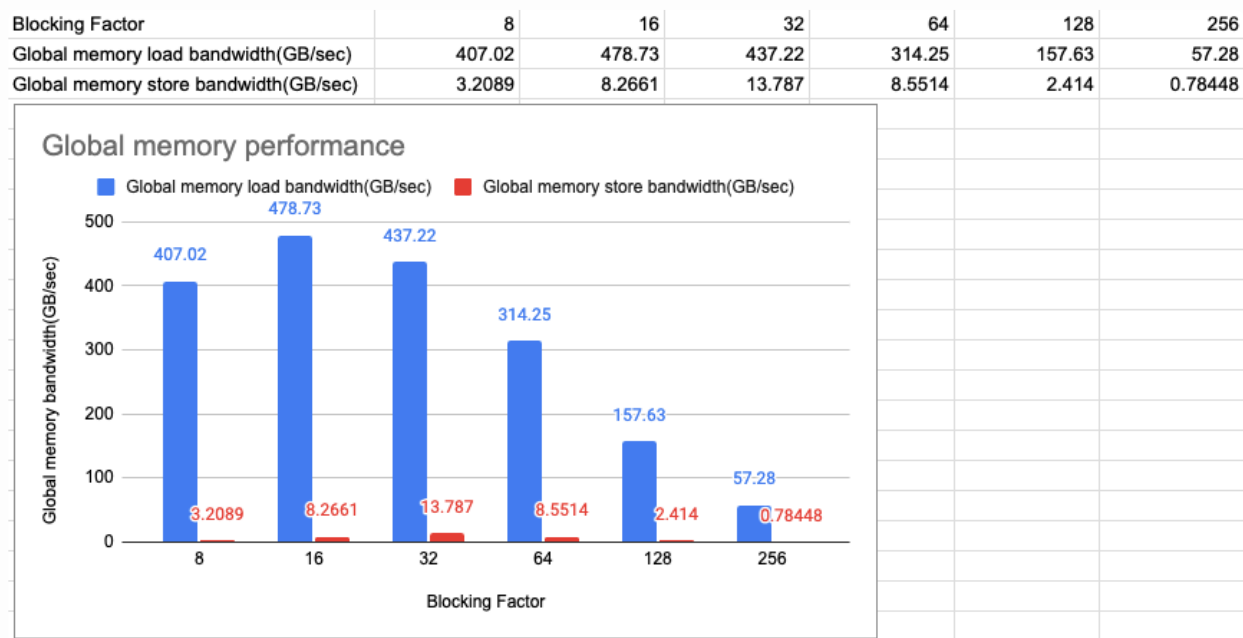


▲不同blocking factor影響GPU warp的使用效率

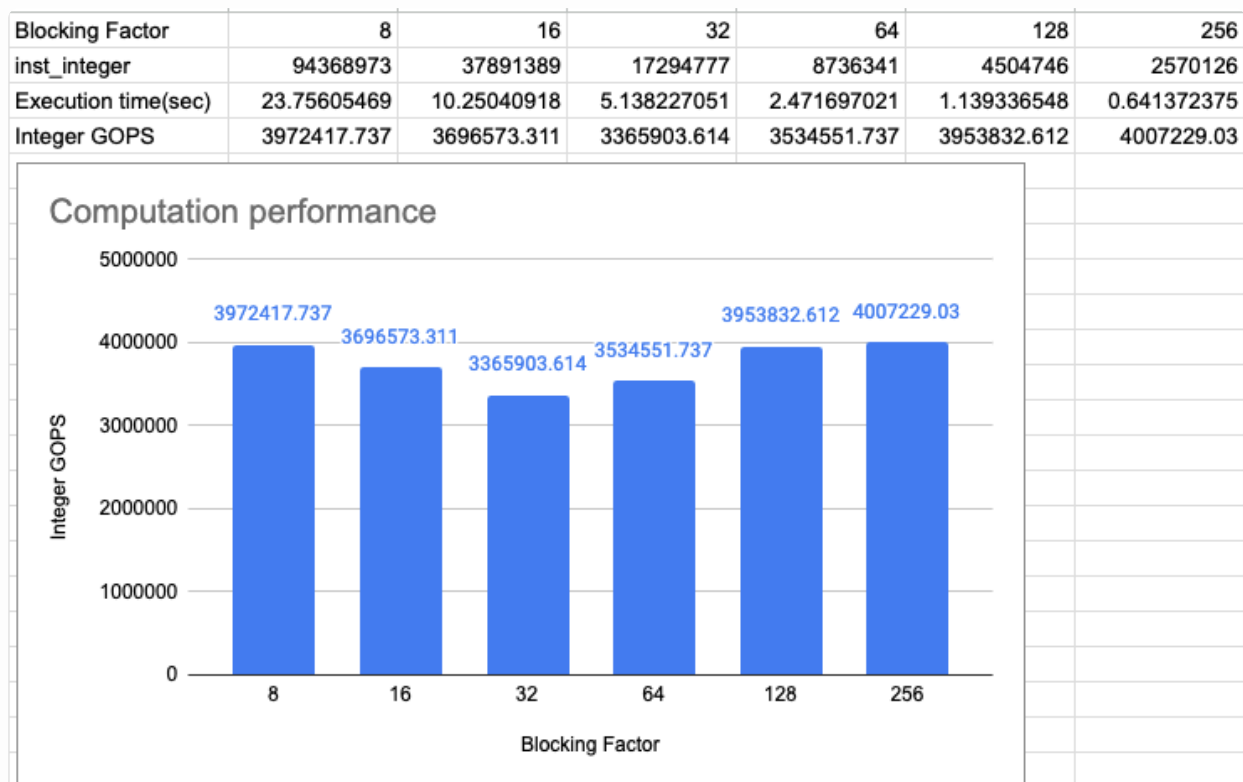


▲不同blocking factor影响GPU core的佔用率和Streaming Multiprocessor的使用效率

Experiment & Analysis



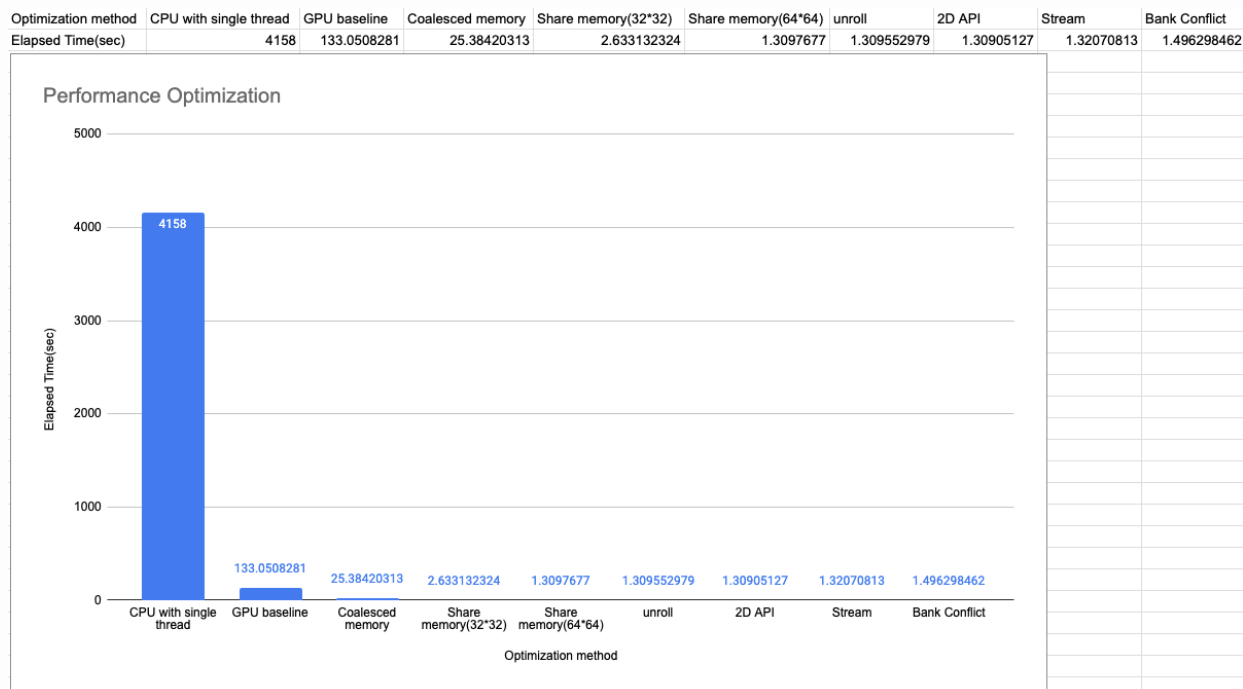
▲不同blocking factor影响GPU global memory load和store操作的bandwidth



▲不同blocking factor影響GPU 每秒整數計算量

Optimization

用p11k1測資當作input資料，測試每多增加一種優化方法，能加速多少



由於CPU會超時，所以實驗時是把(一輪的時間 X round)當作結果

CPU baseline是完全沒有優化的方法

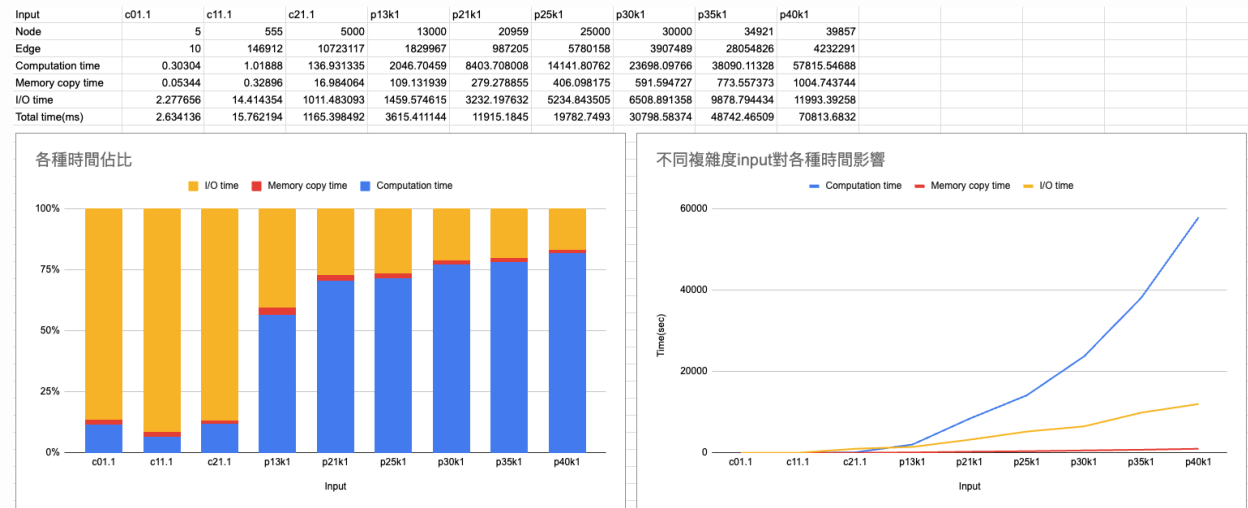
GPU baseline是把每個pixel都丟給一個GPU thread計算，所以可以快很多倍

後面實作unroll、stream、bank conflict方法好像對整體影響不大，甚至有可能讓整體效能變差

Time Distribution

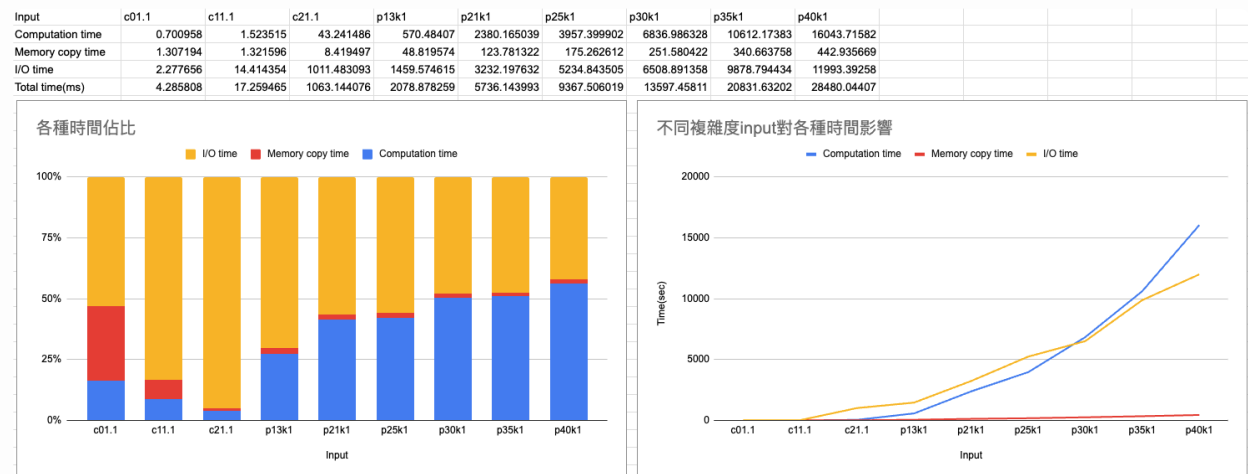
根據不同的scale選擇了不同的測資當作input

單GPU實驗



在計算工作複雜度較小時，I/O幾乎等於整體運算時間。但隨著vertex和edge的數量上升，計算量變大，計算才逐漸變成程式的bottleneck。

單GPU跑AMD實驗



一樣的程式轉成使用AMD的API，並跑在AMD的GPU上，運算速度竟然快了好多。當測到p40k1這筆測資時，I/O時間居然和計算時間沒差多少。而且在做GPU memory copy的時候，還發現nvidia的copy from host to device和device to host的時間差距其實蠻大的。AMD的API對這兩個操作所花的時間幾乎是一樣的。

雙GPU實驗



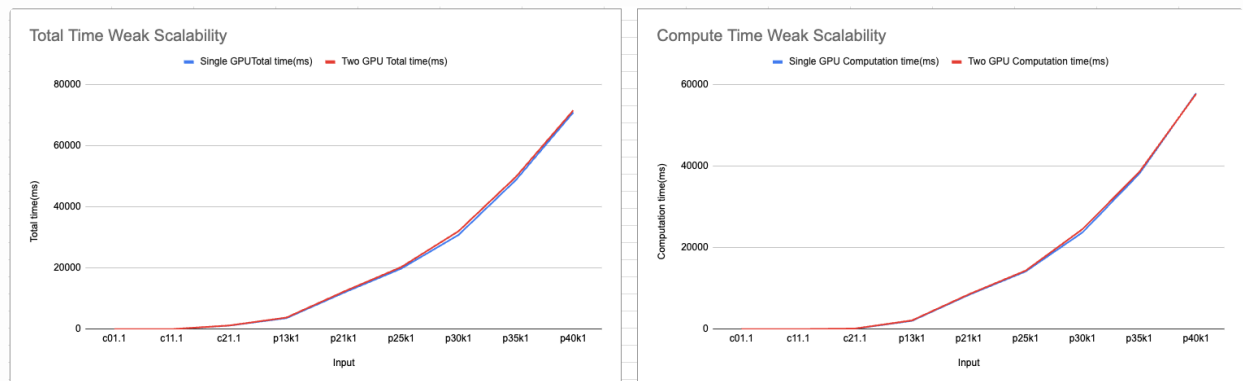
多GPU和單GPU的狀況一樣，都是計算工作的scale決定計算時間和I/O的佔比。

雙GPU跑AMD實驗

```
pp24s086@apollo-login:~/hw3-3$ srun -N1 -n1 -c2 -p amd --gres=gpu:2 hw3-3-amd testcases/c01.1 c01.1.0
ut
srun: error: apollo-amd-vm: task 0: Exited with exit code 127
hw3-3-amd: error while loading shared libraries: libomp.so: cannot open shared object file: No such file or directory
```

雙GPU跑編譯好的AMD code不知道為什麼會遇到不能用OpenMP的問題，所以就沒有測了。

Weak Scalability



對於單GPU和雙GPU而言差距不算太大，但因為實作的關係，多GPU在運作上還比單GPU慢一點點。如果只看計算不看I/O，結果也是一樣。

Conclusion

這次作業有讓我嘗試了很多上課教到的GPU優化技巧。經過不同的實作，我也發現不同的情境下並不是每個優化技巧都是用，有些技巧如stream會需要花時間創建一個stream，若stream只用沒幾次就不用了，會很浪費。又或是2D memory copy API會需要先pin host memory再對齊，才能讓API加速。如果要copy的資料很小且很雜，那麼pin memory加上對齊的時間就會搞垮整個效率，因此要慎重定使用的優化技巧。多GPU在實作後才發現，以目前所學不太有優化的空間，基本上沒有優化還可能比單GPU還糟。

