

1. 병렬화방식

우선 **A,B,C**의 틀은 모든 노드에 있고, 실제 데이터(**A,B**)들은 노드0에만 있는 상태다. 노드 4개, 그리고 각 노드에 달린 GPU에 이제 행렬곱의 일을 분배할 것인데, 일단 **B**행렬은 모든 노드에 **broadcast**하여 전부 보낸다.

이후 행렬 **A**의 행을 노드 개수만큼 나누어, 즉, $M/4 * K$ 만큼의 일을 각 노드에게 분배한다(**scatter**). 다만 4로 나누어 떨어지지 않는 **M**에 대해서는 마지막 노드에게 4로 나눈 나머지만큼 할당하여 분배한다.

각 노드에 대해서는, 분배받은 **M/4**개(except for last node) 행에 대해 또 달린 GPU개수만큼 나누어(4개) 각 GPU가 일을 하도록 분배했다.

그렇게 각 노드별, 각 GPU들이 분배받은 작은 **sub**행렬곱을 실행한 결과를 마지막에 노드0에 보내서(**Gather**) **C**를 완성한다.

2. 뼈대코드 설명

- **matmul_initialize():** MPI 초기세팅(**community**설정, **node**개수, **gpu**디바이스 개수확인, **gpu**디바이스별로 **global memory**에 값을 받아올 공간 할당)
 - **CudaGetDeviceCount(&num_devices):** num_device에 노드당 달려있는 (cuda가 지원하는) **gpu**디바이스개수 받기(개수조회)
 - **CudaGetDeviceProperties(&prop,i):** 자신의 GPU디바이스(**compute device**) 정보 조회
 - **cudaSetDevice(i):** **gpu**디바이스가 **gpu execution**에 활용될 수 있도록 세팅
 - **cudaMalloc():** **gpu**디바이스의 **global memory(DRAM)**에 메모리공간 할당
- **matmul():** 행렬곱을 노드 간, 노드 내 **gpu**별로 일을 분배해 병렬적으로 실행. **MPI**를 이용해 **A,B**의 데이터들을 노드에 보내고 **kernel** 함수를 런칭한다.
 - **cudaMemcpy(void* dist, void* src, size_t count, datatype):** GPU메모리에서 CPU메모리로, 또는 반대방향으로 메모리를 복사하거나 전송할 때 쓴다.
 - **cudaDeviceSynchronize():** 한 디바이스in **gpu**가 기존에 **requested tasks**를 모두 끝낼 때까지 **block**한다.
- **matmul_finalize():**
 - **matmul_initialize()**에서 메모리 할당했던(**cudaMalloc()**으로) 메모리들을 모두 **free**해준다.
 - **cudaFree():** 사용이 끝난 메모리 해제

3. 최적화 방식 및 실험 결과

- **shared memory** (in SM in **gpu** device)에 **global memory**로부터 갖고온 행렬의 일부를 **sub matrix**([tilesize][tilesize] 크기) 로 담아두어 **global memory**에 대한 접근을 줄인다. (**tiling + reuse**비율 높이기)
- 실험1. **shared memory**에 행렬을 담아올 때, 행 기준으로 **threadidX.x**부터 받아오는 것 vs **threadidX.y**부터 받아오는 것
 - 후자의 경우 **threadidX.y**가 **y**축부분에 해당하므로 **row-based**로 접근하게 된다.
 - 전자보다 후자가 빠르다.
 - 성능평가로 준 코드에서 전자: 약 10300GFLOPS, 후자: 약 12000GFLOPS

- 실험2. tile_size=16 vs tile_size = 32
 - shared memory로 담아올 때의 크기를 32와 16으로 두고 실험해봤다.
 - 32가 16보다 빠르다. (아마, 16인 경우보다 32인 경우가 global memory에 접근하는 것이 더 많아서 인듯 하다.
 - ./run.sh -v -n 1 2227 3493 2999 기준, 전자: 약 2286 GFLOPS 후자: 약 2500GFLOPS
- 4. OpenCL vs CUDA
 - OpenCL에 비해 CUDA는 호스트코드, 커널코드를 따로 나눠서 작성하지 않아도 같이 컴파일이 되어서 좋다.
 - CUDA에 비해 OpenCL은 Constant 메모리의 동적할당이 가능하단 장점이 있다.