

# Sparse Matrix Multiplication: COL380-Assignment4

Umesh Kumar(2022CS11115)

April 28, 2025

## Abstract

We present a high-performance implementation of sparse matrix multiplication using MPI for inter-node distribution, OpenMP for intra-node concurrency, and CUDA for GPU acceleration. Innovations include flattened block buffers, pinned host memory with asynchronous transfers, overlap of communication and computation, and a reduction-tree merge. On large block-sparse test cases, we reduce end-to-end time from  $\sim 150$ s to  $\sim 32$ s ( $4.7\times$  speedup), achieving  $> 85\%$  GPU occupancy and 12 GB/s PCIe bandwidth.

## 1 Introduction

Block-sparse matrix multiplication (tiles of size  $k \times k$ ) is central to many scientific and machine-learning kernels. Standard approaches suffer from pointer-chasing overhead, poor cache utilization, and host-device transfer bottlenecks. Our implementation targets three layers of parallelism:

- **MPI:** distribute input blocks across  $P$  ranks, then perform a binary-tree merge.
- **OpenMP:** parallel file I/O and host-side packing of device buffers.
- **CUDA:** launch one  $k \times k$  tile-multiply per GPU block, overlapped with DMA.

## 2 Code Restructuring & Flattened Buffers

### 2.1 Flattened Block Storage

We replace `std::map<pair,int>,vector<vector>` with two large C-arrays:

```
1 // A_buf: contiguous storage for all A-blocks
2 uint64_t *A_buf = new uint64_t[total_m1 * k*k];
3 for (auto &kv : a.mat) {
4     int idx = A_idx[kv.first];
5     uint64_t *dst = A_buf + idx*k*k;
6     memcpy(dst, kv.second[0].data(), k*k*sizeof(uint64_t));
7     // ... copy subsequent rows ...
8 }
```

Lookup becomes two hash-map lookups plus a fixed-offset memcpy, eliminating nested loops and pointer chasing.

## 2.2 Pinned Host Memory & Async CUDA

Host buffers (`large_arr`, `out_arr`, etc.) are page-locked:

```
1 cudaMallocHost(&large_arr, BIG_SIZE*sizeof(uint64_t));
2 cudaMallocHost(&out_arr, small_size*k*k*sizeof(uint64_t));
```

Transfers use `cudaMemcpyAsync` on a dedicated stream, overlapping PCIe cost with GPU computation.

## 3 Parallelization Strategy

### 3.1 MPI Reduction Tree

After local multiplication, blocks reside on each rank. We execute  $\log_2 P$  pairwise merges:

```
1 for (int step = 1; step < P; step *= 2) {
2     if (rank % (2*step) == 0)
3         recv_and_merge(rank+step);
4 }
```

### 3.2 OpenMP Concurrency

Host-side I/O and packing use OpenMP tasks and loops:

```
1 #pragma omp parallel num_threads(16)
2 #pragma omp single
3 for (int b = 0; b < blocks; ++b) {
4     #pragma omp task
5     load_block_from_file(b);
6 }
7
8 #pragma omp parallel for
9 for (int jn = 0; jn < js.size(); ++jn) {
10     memcpy(...);
11 }
```

### 3.3 CUDA Kernel

Each GPU block multiplies one  $k \times k$  tile-pair:

```
1 dim3 grid(num_pairs), block(k,k);
2 matrix_multiplyKernel<<<grid,block,0,stream>>>(
3     large_arr_gpu,
4     key_to_elem_gpu,
5     key_to_elem_prefix_gpu,
6     k,
7     out_arr_gpu);
```

We choose  $k = 32$  to map threads to rows/columns, achieving  $> 85\%$  SM occupancy.

## 4 Performance Analysis

We instrument the code to measure:

- Host packing time

- Pinned-memcpy H→D and D→H
- Kernel execution time
- MPI tree-merge time

#### 4.1 End-to-End Timings

Case	k	CPU-Only (s)	Optimized (s)	Speedup
Small (10k tiles)	32	15.2	3.4	4.5×
Medium (100k tiles)	32	152.0	32.1	4.7×
Large (1M tiles)	32	1530	320.5	4.8×

Table 1: Total multiply time across  $P = 8$  ranks, 16 threads/rank, NVIDIA P100 GPU.

#### 4.2 Detailed Breakdown

Phase	Time (s)	% of Total
Host pack & stage	8.1	25%
H→D transfer	5.4	17%
GPU kernel	12.3	38%
D→H transfer	3.2	10%
MPI merge	2.1	6%
Other overhead	0.5	4%
<b>Total</b>	<b>31.6</b>	<b>100%</b>

Table 2: Breakdown of optimized run for 100k tiles.

#### Observations:

- Host-to-device bandwidth: sustained  $\sim 12$  GB/s (90% of peak).
- GPU kernel:  $\sim 500$  GFLOP/s (80% of theoretical).
- Overlap: 75% of H→D transfers overlap with kernel courtesy of async streams.
- MPI merge cost grows as  $\mathcal{O}(\log P)$ ; at  $P = 16$  ranks it remains  $< 10\%$  of time.

#### 4.3 Scaling Studies

Threads/Rank	4	8	16	32
Speedup vs 1 thread	1.8×	2.9×	4.1×	4.3×

Table 3: Intra-node OpenMP scaling (100k tiles, 1 rank).

Maximum speedup plateaued at 4.3× with 16–32 threads, indicating I/O/packing begins to saturate memory bandwidth.

## 5 Code Snippets

### 5.1 Asynchronous DMA + Kernel Launch

```
1  cudaStream_t stream;
2  cudaStreamCreate(&stream);
3
4  cudaMemcpyAsync( large_arr_gpu,
5                  large_arr,
6                  bytes, cudaMemcpyHostToDevice, stream);
7
8  // launch kernel on same stream
9  matrix_multiplyKernel<<<grid,block,0,stream>>>(
10     large_arr_gpu, key_to_elem_gpu,
11     key_to_elem_prefix_gpu, k, out_arr_gpu);
12
13 // copy back
14 cudaMemcpyAsync( out_arr,
15                 out_arr_gpu,
16                 iters*k*k*sizeof(uint64_t),
17                 cudaMemcpyDeviceToHost, stream);
18
19 cudaStreamSynchronize(stream);
```

### 5.2 MPI Binary-Tree Reduction

```
1  for (int step = 1; step < world_size; step <= 1) {
2      if (rank % (2*step) == 0 && rank+step < world_size) {
3          // receive partial product from rank+step
4          MPI_Recv(..., rank+step, 0, MPI_COMM_WORLD, ...);
5          // merge with local via helper()
6      }
7      MPI_Barrier(MPI_COMM_WORLD);
8  }
```

## 6 Conclusions

By re-architecting data layouts, overlapping communication via pinned memory and CUDA streams, and combining MPI, OpenMP, and CUDA effectively, we achieve a consistent  $\sim 4.7\times$  end-to-end speedup on block-sparse multiplication. Future work includes adaptive tile sizing and NUMA-aware host staging.