

Sparse Matrix-Vector Multiplication: Performance Analysis and Optimization

Ahmet YOLDAS
AYBU - Parallel Algorithms
Final Project Report

December 2025

Abstract

This report presents a comprehensive performance analysis of Sparse Matrix-Vector Multiplication (SpMV) implementations using various optimization techniques. We implemented and evaluated five methods: CSR Serial (baseline), CSR Parallel, BCSR Parallel, CSR+Bucket Parallel, and BCSR+Bucket Parallel. Our experimental results on 2000×2000 random sparse matrices (5% density) demonstrate that simple CSR parallelization achieves 14.37 GFlop/s with $4.84 \times$ speedup, significantly outperforming complex methods like BCSR (1.09 GFlop/s). We developed a novel adaptive bucket sizing algorithm that improved CSR+Bucket performance by $7.5 \times$ (from 1.62 to 12.18 GFlop/s), achieving 85% of CSR Parallel performance. Our analysis reveals that matrix characteristics critically affect method selection: random sparse matrices favor simple CSR parallelization, while structured sparse matrices benefit from blocking techniques. The final implementation is production-ready with modular design, achieving 221% Roofline efficiency.

1 Introduction

Sparse Matrix-Vector Multiplication (SpMV) is a fundamental operation in scientific computing, appearing in iterative solvers, graph algorithms, and machine learning applications. The performance of SpMV directly impacts the overall efficiency of these applications.

1.1 Problem Statement

SpMV computes $\mathbf{y} = A\mathbf{x}$ where A is a sparse $n \times n$ matrix and \mathbf{x} is a dense vector. The key challenges are:

- Irregular memory access patterns
- Low arithmetic intensity (0.125 flops/byte for double precision)

- Memory bandwidth bottleneck
- Load balancing in parallel implementations

1.2 Project Objectives

1. Implement and compare multiple SpMV methods
2. Achieve $> 2\times$ speedup over serial baseline
3. Achieve > 3 GFlop/s performance
4. Analyze matrix characteristics impact on performance
5. Develop production-ready modular code

2 Methods Implemented

We implemented five SpMV methods in a modular design, with each method in separate source files for maintainability and extensibility.

2.1 Method 1: CSR Serial (Baseline)

The Compressed Sparse Row (CSR) format stores only non-zero elements with three arrays:

- `values[]`: Non-zero values
- `col_idx[]`: Column indices
- `row_ptr[]`: Row pointers

Implementation (csr_serial.c):

```

1 void spmv_csr_serial(const CSR_Matrix *A,
2                       const double *x, double *y) {
3     for (int i = 0; i < A->rows; i++) {
4         double sum = 0.0;
5         for (int k = A->row_ptr[i];
6              k < A->row_ptr[i+1]; k++) {
7             sum += A->values[k] * x[A->col_idx[k]];
8         }
9         y[i] = sum;
10    }
11 }
```

Performance: 2.97 GFlop/s (baseline)

2.2 Method 2: CSR Parallel (Final Method)

Parallelizes the outer loop using OpenMP with dynamic scheduling for load balancing.

Implementation (csr_parallel.c):

```
1 void spmv_csr_parallel(const CSR_Matrix *A,
2                         const double *x, double *y) {
3 #pragma omp parallel for schedule(dynamic, 64)
4     for (int i = 0; i < A->rows; i++) {
5         double sum = 0.0;
6         for (int k = A->row_ptr[i];
7              k < A->row_ptr[i+1]; k++) {
8             sum += A->values[k] * x[A->col_idx[k]];
9         }
10        y[i] = sum;
11    }
12 }
```

Key optimizations:

- OpenMP parallel for with 8 threads
- Dynamic scheduling (chunk size = 64)
- Load balancing for irregular row lengths

Performance: 14.37 GFlop/s (4.84 \times speedup)

2.3 Method 3: BCSR Parallel

Block Compressed Sparse Row (BCSR) uses 4×4 register blocking with OpenMP parallelization.

Key features:

- 4×4 dense blocks
- Register reuse
- Loop unrolling
- OpenMP parallelization

Performance: 1.09 GFlop/s (0.37 \times speedup)

Storage overhead: 11.2 \times (200K non-zeros \rightarrow 2.2M storage)

2.4 Method 4: CSR+Bucket Parallel (Optimized)

Implements adaptive bucket partitioning with cache-aware sizing.

Adaptive bucket algorithm:

```
1 int num_threads = omp_get_max_threads();
2 int min_buckets = num_threads * 4; // 32 for 8 threads
3
4 int bucket_size = A->rows / min_buckets;
5 if (bucket_size < 32) bucket_size = 32; // Min
6 if (bucket_size > 512) bucket_size = 512; // Max
7
8 // Result: 62 rows/bucket, 33 buckets for 2000 rows
```

Performance: 12.18 GFlop/s (4.10 \times speedup)

Improvement: 7.5 \times faster than fixed bucket size (1.62 GFlop/s)

2.5 Method 5: BCSR+Bucket Parallel (Hybrid)

Combines 4 \times 4 blocking with adaptive bucket partitioning.

Performance: 2.13 GFlop/s (0.72 \times speedup)

Improvement: 4 \times faster than fixed bucket size (0.53 GFlop/s)

3 Experimental Setup

3.1 Hardware Configuration

- CPU: Intel Core (details from system)
- Threads: 8 OpenMP threads
- Memory: DDR4
- Compiler: GCC 11.4.0
- Optimization: -O3 -march=native -fopenmp

3.2 Test Matrix

- Size: 2000 \times 2000
- Density: 5% (target)
- Actual non-zeros: 200,390
- Actual density: 5.01%

- Pattern: Random sparse
- Generation: Pseudo-random (seed=42)

3.3 Methodology

1. Generate random sparse matrix (CSR format)
2. Convert to BCSR for blocking methods
3. Warm-up run for cache
4. Timed run (single SpMV operation)
5. Verify correctness against serial baseline
6. Compute GFlop/s: $\frac{2 \times \text{nnz}}{\text{time} \times 10^9}$

4 Results

4.1 Performance Summary

Table 1: Performance comparison of all methods

| Method | Time (ms) | GFlop/s | Speedup |
|-------------------------|--------------|--------------|--------------|
| CSR Serial | 0.135 | 2.97 | 1.00× |
| CSR Parallel | 0.028 | 14.37 | 4.84× |
| BCSR Parallel | 0.367 | 1.09 | 0.37× |
| CSR+Bucket (Optimized) | 0.033 | 12.18 | 4.10× |
| BCSR+Bucket (Optimized) | 0.188 | 2.13 | 0.72× |

4.2 Roofline Analysis

Roofline model parameters:

- Peak bandwidth: 52 GB/s
- Arithmetic intensity: 0.125 flops/byte
- Roofline peak: $52 \times 0.125 = 6.5$ GFlop/s

CSR Parallel efficiency: $\frac{14.37}{6.5} = 221\%$

The 221% efficiency indicates our implementation exceeds the conservative Roofline estimate, likely due to:

- CPU frequency boost under load

- Cache reuse (x vector)
- Optimized memory access patterns

4.3 Bucket Optimization Impact

Table 2: Impact of adaptive bucket sizing

| Method | Before | After | Improvement |
|-------------|-----------------------------|-------------------------------|-------------|
| CSR+Bucket | 1.62 GFlop/s (2 buckets) | 12.18 GFlop/s (33 buckets) | 7.5× |
| BCSR+Bucket | 0.53 GFlop/s (2 buckets) | 2.13 GFlop/s (34 buckets) | 4.0× |

5 Analysis

5.1 Why CSR Parallel Won

CSR Parallel achieved the best performance (14.37 GFlop/s) for several reasons:

5.1.1 Simplicity and Efficiency

- Minimal overhead (16 lines of code)
- Direct parallelization of outer loop
- No storage overhead
- Optimal for OpenMP

5.1.2 Effective Load Balancing

Dynamic scheduling (chunk=64) handles irregular row lengths:

- 2000 rows / 64 = 31 chunks
- 31 chunks for 8 threads = 3.9 chunks/thread
- Good granularity for load distribution

5.1.3 Memory Access Patterns

- Sequential access to `values[]` and `col_idx[]`
- Random access to `x[]` (cached across rows)
- Sequential write to `y[]` (no conflicts)

5.1.4 Cache Utilization

- Vector x reused across rows (cache hits)
- Small working set per thread
- No cache line conflicts (independent rows)

5.2 Why BCSR Failed

BCSR Parallel achieved only 1.09 GFlop/s ($0.37\times$ speedup) due to:

5.2.1 Massive Storage Overhead

- Original non-zeros: 200,390
- BCSR storage: $140,132 \text{ blocks} \times 16 = 2,242,112 \text{ elements}$
- Overhead: $\frac{2,242,112}{200,390} = 11.2\times$
- Most blocks are sparse (wasted storage and computation)

5.2.2 Poor Match for Random Sparse

Random sparse matrices have scattered non-zeros:

- 4×4 blocks mostly contain zeros
- Blocking efficiency: $\frac{200,390}{2,242,112} = 8.9\%$
- 91.1% of operations on zeros

5.2.3 Cache Pollution

- $11\times$ more data loaded into cache
- Cache capacity wasted on zeros
- Reduced effective cache size

5.2.4 When BCSR Would Work

BCSR excels on structured sparse matrices:

- Finite element matrices (FEM)
- Banded matrices
- Structured stencils

- Block diagonal systems

Example: 7-point stencil in 3D would have dense 4×4 blocks, achieving 90%+ efficiency.

5.3 Adaptive Bucket Optimization

5.3.1 The Problem

Fixed bucket size (1024 rows) created poor parallelism:

$$\text{Buckets} = \frac{2000}{1024} \approx 2 \text{ buckets} \quad (1)$$

With 8 threads and 2 buckets:

- Only 2 threads active
- 6 threads idle (75% waste)
- No load balancing
- Poor thread utilization

5.3.2 Adaptive Solution

Our algorithm ensures sufficient buckets:

$$\text{min_buckets} = \text{threads} \times 4 = 32 \quad (2)$$

$$\text{bucket_size} = \frac{\text{rows}}{\text{min_buckets}} = \frac{2000}{32} = 62.5 \approx 64 \quad (3)$$

Result: 33 buckets for 8 threads (4.1 buckets/thread)

5.3.3 Performance Impact

- All threads utilized
- Dynamic scheduling effective
- Excellent load balancing
- Cache-friendly (64 rows = 512 bytes)

CSR+Bucket improvement: 1.62 → 12.18 GFlop/s (7.5×)

5.3.4 Why Not Better Than CSR Parallel?

CSR+Bucket (12.18 GFlop/s) still lags CSR Parallel (14.37 GFlop/s):

- Bucket overhead: Division and boundary checks
- CSR Parallel: Direct row parallelization (no overhead)
- Gap: 2.2 GFlop/s (15%)

But CSR+Bucket scales better for very large matrices (future work).

5.4 BCSR+Bucket Hybrid Analysis

BCSR+Bucket combines blocking and buckets:

- Adaptive buckets improved parallelism: $0.53 \rightarrow 2.13$ GFlop/s ($4\times$)
- Still limited by BCSR overhead ($11\times$ storage)
- Better than BCSR alone (1.09 GFlop/s)
- Shows double optimization has limits

Key insight: Base overhead (BCSR storage) dominates, preventing hybrid from competing with simple CSR methods.

6 Matrix Characteristics Impact

6.1 Random Sparse Matrices

Our test case: 2000×2000 , 5% density, random

Winner: CSR Parallel (14.37 GFlop/s)

Why:

- Scattered non-zeros
- BCSR blocks mostly empty (8.9% efficiency)
- Simple parallelization optimal
- No structure to exploit

6.2 Structured Sparse Matrices

Examples: FEM, stencils, band matrices

Expected winner: BCSR Parallel

Why:

- Dense blocks (90%+ efficiency)
- Register reuse effective
- SIMD vectorization
- Reduced storage overhead

Example performance (hypothetical 7-point stencil):

- BCSR Parallel: ~8-10 GFlop/s
- CSR Parallel: ~6-7 GFlop/s

6.3 Large Sparse Matrices

Example: $100,000 \times 100,000$, 0.01% density

Expected winner: CSR+Bucket Parallel

Why:

- Bucket count: $\frac{100,000}{512} \approx 200$ buckets
- Excellent load balancing
- Better cache locality (bucket-sized chunks)
- Scales well to many cores

6.4 Small Dense Matrices

Example: 500×500 , 80% density

Best approach: Dense BLAS (not sparse methods)

Why:

- Sparse overhead not justified
- Dense methods achieve near-peak
- Cache blocking effective

Table 3: Method selection guide

| Matrix Type | Best Method | Expected Perf. |
|-----------------------|---------------|----------------|
| Random sparse (small) | CSR Parallel | 10-15 GFlop/s |
| Random sparse (large) | CSR+Bucket | 8-12 GFlop/s |
| Structured (FEM) | BCSR Parallel | 8-10 GFlop/s |
| Banded | BCSR Parallel | 6-9 GFlop/s |
| Very sparse (<1%) | CSR Parallel | 5-8 GFlop/s |
| Dense-like (>50%) | Dense BLAS | 30+ GFlop/s |

6.5 Decision Matrix

7 Progressive Development

Our implementation evolved through multiple iterations:

7.1 Checkpoint-1: Initial Implementation

- CSR Serial: 1.72 GFlop/s
- CSR Parallel: 5.20 GFlop/s ($3.02\times$)
- Basic OpenMP parallelization
- Static scheduling

7.2 Checkpoint-2: Method Comparison

- Added BCSR Parallel: 1.11 GFlop/s
- Identified BCSR overhead issue
- Improved CSR Parallel: 7.86 GFlop/s ($4.57\times$)
- Dynamic scheduling optimization

7.3 Modular Design

- Separated each method into own files
- 13 source files (professional structure)
- Easy to modify and extend
- Production-ready organization

7.4 Bucket Methods

- Initial fixed bucket: 1.62 GFlop/s (poor)
- Problem identified: Only 2 buckets
- Adaptive algorithm developed
- Final CSR+Bucket: 12.18 GFlop/s ($7.5\times$ improvement)

7.5 Final Optimization

- CSR Parallel: 14.37 GFlop/s ($4.84\times$)
- All methods optimized
- Comprehensive analysis
- Publication-quality results

8 Conclusions

8.1 Key Findings

8.1.1 Performance Achievement

- **Final method:** CSR Parallel
- **Performance:** 14.37 GFlop/s
- **Speedup:** $4.84\times$ vs serial baseline
- **Efficiency:** 221% Roofline
- **Exceeded targets:** $4.84\times$ vs $2\times$ goal (242% of target)

8.1.2 Algorithm Selection

Simple CSR parallelization proved optimal for random sparse matrices:

- Minimal code (16 lines)
- No storage overhead
- Effective load balancing
- Production-ready

8.1.3 Adaptive Bucket Innovation

Our adaptive bucket sizing algorithm:

- $7.5\times$ improvement for CSR+Bucket
- Matrix and thread-aware
- Ensures optimal parallelism
- Original contribution

8.1.4 Matrix Characteristics Matter

Performance critically depends on matrix structure:

- Random sparse: CSR Parallel best
- Structured sparse: BCSR competitive
- Large matrices: Bucket methods scale better
- No universal best method

8.1.5 Complexity vs Performance

Complex optimizations don't always win:

- CSR Parallel (16 lines): 14.37 GFlop/s
- BCSR+Bucket (66 lines): 2.13 GFlop/s
- Simplicity often beats complexity
- Overhead analysis essential

8.2 Scientific Contributions

1. **Adaptive bucket sizing algorithm:** Ensures $4\times$ buckets per thread, adapting to matrix size
2. **BCSR overhead quantification:** $11.2\times$ storage overhead on random sparse, 8.9% efficiency
3. **Hybrid method analysis:** BCSR+Bucket shows double optimization limits
4. **Matrix-method mapping:** Clear guidelines for method selection
5. **Production code:** Modular, documented, tested implementation

8.3 Future Work

- Test on larger matrices (100K+ rows)
- Evaluate on structured sparse (FEM, stencils)
- GPU implementation
- Auto-tuning for matrix characteristics
- Integration with sparse solvers

8.4 Final Recommendation

For production use on random sparse matrices:

- Use: CSR Parallel
- Performance: 14.37 GFlop/s (proven)
- Simplicity: 16 lines, easy to maintain
- Portability: Standard OpenMP
- Robustness: Tested and verified

Source code available at: `spmv_optimized_final/`

References

- [1] Azad, A., & Buluç, A. (2016). A work-efficient parallel sparse matrix-vector multiplication algorithm. *IEEE International Parallel and Distributed Processing Symposium*.
- [2] Vuduc, R., Demmel, J. W., & Yelick, K. A. (2002). OSKI: A library of automatically tuned sparse matrix kernels. *In Proceedings of SciDAC*.
- [3] Im, E. J., & Yelick, K. (2004). Optimizing sparse matrix computations for register reuse in SPARSITY. *International Conference on Computational Science*.
- [4] Williams, S., Waterman, A., & Patterson, D. (2009). Roofline: An insightful visual performance model for multicore architectures. *Communications of the ACM*, 52(4), 65-76.

A Source Code Summary

A.1 Modular Structure

```
spmv_optimized_final/
common.h/c           (145 lines) - Matrix utilities
csr_serial.c/h       (13 lines)  - Method 1
csr_parallel.c/h     (16 lines)  - Method 2 (FINAL)
bcsr_parallel.c/h    (45 lines)  - Method 3
bucket_parallel.c/h  (55 lines)  - Method 4 (adaptive)
bcsr_bucket_parallel.c/h (66 lines) - Method 5 (hybrid)
benchmark.c          (275 lines) - Testing framework
plot_results.py      - Visualization
Makefile              - Build system
run_all.sh           - Automation
```

A.2 Compilation

```
gcc -O3 -march=native -fopenmp *.c -o benchmark -lm
./benchmark 2000 0.05 8
python3 plot_results.py
```

B Complete Results Table

Table 4: Detailed performance metrics

| Method | Time (ms) | GFlop/s | Speedup | Efficiency (%) | Code (lines) |
|---------------|--------------|---------|---------|-------------------|-----------------|
| CSR Serial | 0.135 | 2.97 | 1.00× | - | 13 |
| CSR Parallel | 0.028 | 14.37 | 4.84× | 221% | 16 |
| BCSR Parallel | 0.367 | 1.09 | 0.37× | 17% | 45 |
| CSR+Bucket | 0.033 | 12.18 | 4.10× | 187% | 55 |
| BCSR+Bucket | 0.188 | 2.13 | 0.72× | 33% | 66 |