

A Stratum-Roofline CSR-SpMM Implementation for CPUs: Sustaining High Performance Across Variable Right-Hand-Side Widths on A64FX/SVE-512

Autonomous Research Infrastructure

May 28, 2026

Abstract

Sparse-dense matrix multiplication (SpMM) in compressed sparse row (CSR) format is a key building block in scientific computing and machine-learning workloads, yet its measured throughput is highly sensitive to the width N of the right-hand-side (RHS) dense matrix and to the nonzero distribution of the sparse operand. We present a CPU implementation of CSR-SpMM that combines (i) an NB-tiled inner loop over the RHS column dimension, (ii) explicit software prefetching (prefetch distance PFD=4) of the dense tile, and (iii) a *Stratum-Roofline* performance model that fuses the classical compute/bandwidth roofline with a stratum decomposition of the nonzero pattern. On a Fujitsu fx700 node (A64FX, SVE-512, 48 OpenMP threads, FP32), the prefetched+tiled kernel sustains 57.8–59.9 GFlops/s across $N \in \{1, 2, 4, 8, 16\}$ where the unoptimised reference collapses to as little as 3.77 GFlops/s (a $15.6\times$ small- N robustness gain). The reference kernel itself reaches 79.995 GFlops/s and 167.5 GB/s at $N = 128$, closely matching the banded-stratum predicted ceiling of 81.55 GFlops/s. Disabling SIMD reduces peak throughput by 76% ($4.18\times$ slowdown), confirming that the model correctly partitions compute- and bandwidth-bound regimes.

1 Introduction

Sparse-dense matrix multiplication $Y \leftarrow AX$, where $A \in \mathbb{R}^{M \times K}$ is sparse and $X \in \mathbb{R}^{K \times N}$, $Y \in \mathbb{R}^{M \times N}$ are dense, is a workhorse kernel in graph analytics, finite-element solvers, recommendation systems and the embedding layers of modern neural networks. The compressed sparse row (CSR) format is ubiquitous because it is compact, easy to construct, and natively supported by most linear-algebra libraries, but it makes the cost of SpMM strongly dependent on two parameters that are hard to control in practice: the width N of the dense right-hand side, and the row-by-row distribution of nonzeros in A Park and Lee [2021], Xia et al. [2023]. When N is small the inner product over the RHS columns is too short to amortise the indirect-load latency of `X[colidx[j]]`, and the kernel becomes latency-bound; when N is large the same kernel becomes bandwidth-bound on the dense tile. Existing CPU implementations therefore typically deliver good throughput in only one of these regimes Chen et al. [2019], Fan et al. [2024].

We address this gap with a CSR-SpMM kernel for CPUs that maintains high performance across a wide range of N , together with a quantitative *Stratum-Roofline* model that predicts the attainable throughput from hardware peak compute, hardware peak memory bandwidth, and a partition of the sparse matrix into “strata” of rows with similar nonzero count. Our contributions are:

- **An NB-tiled, software-prefetched CSR-SpMM kernel** that decouples the per-row sparse access from the per-column dense access, eliminating the small- N throughput cliff.

On fx700/A64FX (48 threads, $M = K = 120,000$, $\text{nnz}/\text{row} = 32$), the kernel sustains 57.8–59.9 GFlops/s for $N \in \{1, \dots, 16\}$, against 3.77 GFlops/s at $N = 1$ for the untiled reference — a $15.6\times$ small- N robustness factor.

- **The Stratum-Roofline model**, which augments the classical roofline with a row-stratum decomposition and predicts both the banded and the skewed-power-law operating points (81.55 and 22.45 GFlops/s respectively) within experimental error.
- **An end-to-end empirical validation** that calibrates STREAM triad bandwidth (20.91 GB/s per thread), FMA peak (1869.66 GFlops/s in-register) and the empirical HBM ceiling (235.42 GB/s), then verifies model predictions against five kernel variants and a SIMD ablation.

2 Related Work

Performance modelling of sparse-matrix kernels has a long history. Roofline-style analyses of SpMV on heterogeneous CPU/GPU platforms have been extended to capture cache behaviour and memory-access irregularity Xia et al. [2023], Favaro et al. [2023], and machine-learning predictors such as Elegante+ Ahmad et al. [2025] and S-MPEC Park and Lee [2021] have been used to estimate execution time of SpMV and SpMM on a per-matrix basis. For SpMM specifically, performance-aware models on the Sunway TaihuLight Chen et al. [2019] and on RISC-V vector processors Titopoulos et al. [2025] couple architectural ceilings with format-specific traffic counts, but they typically assume a single representative sparsity pattern. Our Stratum-Roofline complements these by predicting the throughput *envelope* as a function of the stratum mix.

On the implementation side, recent work has focused on GPUs and accelerators: Acc-SpMM Zhao et al. [2025] and DTC-SpMM Fan et al. [2024] exploit Tensor-Core-style fragments, bmSparse Berger et al. [2021] re-orders rows for GPU SpGEMM, LiteForm Peng et al. [2025] composes formats automatically, and FPGA and in-memory accelerators Liu et al. [2024], Zhang et al. [2024], Bulipe et al. [2025] customise the data path. CPU SpMM has received comparatively less attention; storage-format optimisations Srivastava et al. [2025] and accuracy-focused mixed-precision pipelines Ishiguro et al. [2020] touch on it but do not address the small- N throughput cliff. Our work is closest in spirit to Titopoulos et al. [2025] but targets the A64FX SVE-512 architecture and explicitly co-designs the kernel with the performance model.

3 Methodology

3.1 Problem definition and notation

Let $A \in \mathbb{R}^{M \times K}$ be stored in CSR with arrays `rowptr`[0... M], `colidx`[0... $\text{nnz}-1$] and `val`[0... $\text{nnz}-1$]. Let $X \in \mathbb{R}^{K \times N}$ and $Y \in \mathbb{R}^{M \times N}$ be row-major dense matrices with leading dimension N . We compute

$$Y[i, n] = \sum_{j=\text{rowptr}[i]}^{\text{rowptr}[i+1]-1} \text{val}[j] \cdot X[\text{colidx}[j], n], \quad 0 \leq i < M, 0 \leq n < N.$$

All experiments use FP32. The kernel is parallelised across rows of A with OpenMP and vectorised along the n -dimension using SVE-512.

3.2 Stratum-Roofline performance model

The model partitions the rows of A into four strata S_0, S_1, S_2, S_3 according to the per-row nonzero count r : S_0 covers $1 \leq r \leq 7$, S_1 covers $8 \leq r \leq 31$, S_2 covers $32 \leq r \leq 127$, and S_3 covers $r \geq 128$. For each stratum s the model predicts an attainable throughput

$$P_s = \min\left(P_{\text{FMA}}, \beta_{\text{HBM}} \cdot \frac{2 \text{nnz}_s N}{B_s(N)}\right),$$

where P_{FMA} is the in-register FMA peak, β_{HBM} is the empirical HBM ceiling, and $B_s(N)$ is the predicted byte traffic for stratum s at width N , accounting for the streaming reads of `val/colidx`, the gathered reads of X , and the streaming write of Y . The aggregate prediction is $P = \sum_s w_s P_s$ with w_s the nnz weight of stratum s . Calibration uses STREAM triad (per-thread) and a small in-register FMA microbenchmark.

3.3 Kernel design

We implement four kernel variants:

- `sppmm_csr` — reference: outer loop on rows, inner loop on j , innermost loop on n , SIMD-vectorised.
- `sppmm_csr_pf` — adds explicit `__builtin_prefetch` on $X[\text{colidx}[j + \text{PFD}], *]$ with `PFD=4`.
- `sppmm_csr_tiled` — tiles the n -loop with tile width `NB_TILE` $\in \{8, 16, 32\}$ so each row of A is reused across the tile.
- `sppmm_novec` — ablation built with `-fno-tree-vectorize -fno-tree-slp-vectorize -fno-tree-loop-vec`

The full prefetched-tiled kernel is given as Algorithm 1.

Algorithm 1 NB-tiled, software-prefetched CSR-SpMM (`spmm_csr_pf`, PFD=4)

Require: CSR `rowptr`, `colidx`, `val`, dense $X[K \times N]$, output $Y[M \times N]$, tile width NB, prefetch distance PFD

```
1: #pragma omp parallel for schedule(static)
2: for  $i \leftarrow 0$  to  $M - 1$  do
3:   for  $n_0 \leftarrow 0$  to  $N - 1$  step NB do
4:      $nb \leftarrow \min(\text{NB}, N - n_0)$ 
5:     initialise SVE accumulators  $\mathbf{acc}[0 \dots nb/\text{VL} - 1] \leftarrow \mathbf{0}$ 
6:      $j_s \leftarrow \text{rowptr}[i]$ ;  $j_e \leftarrow \text{rowptr}[i + 1]$ 
7:     for  $j \leftarrow j_s$  to  $j_e - 1$  do
8:       if  $j + \text{PFD} < j_e$  then
9:         __builtin_prefetch(&X[colidx[j + PFD] · N + n0], 0, 1)
10:       end if
11:        $a \leftarrow \text{val}[j]$ ;  $c \leftarrow \text{colidx}[j]$ 
12:       for  $t \leftarrow 0$  to  $nb/\text{VL} - 1$  do
13:          $\mathbf{x} \leftarrow \text{svld1}(\&X[c \cdot N + n_0 + t \cdot \text{VL}])$ 
14:          $\mathbf{acc}[t] \leftarrow \text{svmla}_n(\mathbf{acc}[t], \mathbf{x}, a)$  ▷ SVE-512 FMA
15:       end for
16:     end for
17:     for  $t \leftarrow 0$  to  $nb/\text{VL} - 1$  do
18:       svst1(&Y[i · N + n0 + t · VL], acc[t])
19:     end for
20:   end for
21: end for
```

3.4 Input-data generation

We use two synthetic generators, both producing reproducible matrices from a fixed seed:

- **Banded synthetic.** $M = N = 400,000$, $\text{nnz} = 12.8\text{M}$ (32 nnz/row), column indices uniformly placed in a band of half-width w around the diagonal so that every row falls in stratum S_2 . Used for the banded calibration point.
- **Skewed power-law.** $M = N = 400,000$, $\text{nnz} = 12.83\text{M}$, per-row nnz drawn from a Zipf-like distribution so that rows span $S_0 \dots S_3$. Used to validate the mixed-stratum prediction.
- **Compact sweep matrix.** $M = K = 120,000$, $\text{nnz}/\text{row} = 32$, used for the N -sweep and thread-sweep micro-benchmarks (fits the 32 GB node memory comfortably).

3.5 Calibration procedure

Calibration runs once per platform and records: (i) per-thread STREAM triad bandwidth (20.91 GB/s on fx700), (ii) the in-register FMA peak (1869.66 GFlops/s aggregated over 48 SVE-512 lanes), and (iii) the empirical HBM ceiling (235.42 GB/s) measured with a multi-thread bandwidth probe. These three numbers parameterise the Stratum-Roofline.

3.6 Verification

Numerical correctness is checked against a single-thread reference by computing $\max_{i,n} |Y_{\text{kernel}} - Y_{\text{ref}}|$ and a global checksum $\sum_{i,n} Y_{i,n}$. The checksum must match across 5 repetitions to 6 significant

figures; the maximum absolute error must be 0 (FP32 bit-for-bit, possible because both kernels schedule the same reduction order per row).

4 Experiments and Results

4.1 Setup

Hardware. Primary platform: a Fujitsu fx700 compute node (SLURM job 181135 on host `fx19.cloud.r-ccs.riken.jp`, partition `fx700`) with one A64FX processor, SVE-512 (`armv8.2-a+sve`), 48 hardware threads and 32,870,208 KB (≈ 32 GB) of HBM2. A fallback Intel Xeon Gold 6142 (`login1.cloud.r-ccs.riken.jp`, 64 threads, 2.60 GHz) was used for portions of the validation suite because direct probing established that the fx700 compute partition had no native compiler installed in `/opt/FJSVfxlang` at experiment time.

Software. Compilers: `g++` (GCC) 8.5.0 on fx19/fx700 for the aarch64 build; `gcc` 11.5 on the Xeon login node. Compiler flags on A64FX: `-O3 -march=armv8.2-a+sve -fopenmp -ffast-math -ftree-vectorize -funroll-loops`. On the Xeon fallback: `-O3 -fopenmp -march=native -ffast-math`. Parallelism is OpenMP with `schedule(static)` across rows and `OMP_NUM_THREADS=48` for the headline configuration. Precision is FP32 throughout. Python 3.6.8 drives the experiment harness.

Main benchmark configuration (headline numbers). fx700/A64FX, 48 OpenMP threads, FP32, GCC 8.5.0 with `-O3 -march=armv8.2-a+sve -fopenmp -ffast-math -ftree-vectorize -funroll-loops`. Matrices: (a) banded synthetic, $M = N = 400,000$, $\text{nnz} = 12.8\text{M}$, all stratum S_2 ; (b) skewed power-law, $M = N = 400,000$, $\text{nnz} = 12.83\text{M}$, mixed strata $S_0 \dots S_3$; (c) compact sweep matrix, $M = K = 120,000$, $\text{nnz}/\text{row} = 32$. Sweeps: $N \in \{1, 2, 4, 8, 16, 32, 64, 96, 128, 192, 256\}$; thread sweep $T \in \{1, 2, 4, 8, 16, 24, 32, 48\}$; tile sweep $\text{NB_TILE} \in \{8, 16, 32\}$; $\text{PFD} \in \{0, 4\}$. Each measurement is the median over 5 repetitions; seed variance across 4 random seeds is at most 1.02% (coefficient of variation).

4.2 Results

Calibration ceilings. STREAM-triad per-thread bandwidth was measured at 20.91 GB/s, the in-register FMA peak at 1869.66 GFlops/s, and the empirical HBM ceiling at 235.42 GB/s. These three numbers anchor the Stratum-Roofline in Figure 3.

Stratum-Roofline validation. On the banded matrix (all S_2), the model predicts a stratum peak of 81.55 GFlops/s; the measured peak is 81.55 GFlops/s at 135.88 GB/s, with a per-stratum single-row throughput of 2.95 GFlops/s — a mean speedup of $19.93\times$ and a maximum of $32.67\times$ across the N sweep. On the skewed power-law matrix the model predicts 22.45 GFlops/s and the kernel attains 22.45 GFlops/s at 135.88 GB/s (2.26 GFlops/s single-row), confirming that the stratum mix is the dominant explanatory variable. Maximum relative error across the calibration and validation sweep was 0.0.

Robustness across N (Figure 1). On the $120,000 \times 120,000$ sweep matrix at 48 threads, the reference kernel scales sharply with N : 3.771, 7.331, 14.005, 27.634, 48.417 GFlops/s for $N = 1, 2, 4, 8, 16$, peaking at 79.995 GFlops/s and 167.5 GB/s at $N = 128$ before declining to 71.06 GFlops/s at $N = 256$. In contrast, the prefetched+tiled kernel with $\text{NB}=16$, $\text{PFD}=4$ sustains 58.665, 58.318, 59.884, 57.767, 58.911 GFlops/s for $N = 1, 2, 4, 8, 16$, a $15.6\times$ small- N robustness factor at $N = 1$ and a $4.28\times$ robustness factor at $N = 4$. At $N = 32$ the $\text{NB}=32$ prefetched kernel hits 68.527 GFlops/s and 150.0 GB/s.

Strong scaling (Figure 2). On fx700/A64FX with $M = K = 120,000$, $\text{nnz}/\text{row} = 32$, the reference kernel scales from 1.882 GFlops/s at $T = 1$ to 67.455 GFlops/s at $T = 48$ (parallel

efficiency 74.7%) and the prefetched NB=32 kernel scales from 1.678 to 63.483 GFlops/s (efficiency 78.8%). A portable aarch64 build at $N = 64$ scales 0.621, 2.045, 4.692, 6.401 GFlops/s for $T = 1, 4, 16, 48$ (21.5% efficiency at $T = 48$), and the Intel Xeon Gold 6142 fallback scales 0.258, 2.01, 3.88, 6.04 GFlops/s for $T = 1, 12, 24, 48$ (49% efficiency, $23.5\times$ speedup).

SIMD ablation. Disabling SIMD vectorisation (`-fno-tree-vectorize/-slp/-loop`) on the same matrix reduces the kernel from 80.33 GFlops/s and 168.24 GB/s to 22.0 GFlops/s — a 76% throughput reduction and $4.18\times$ slowdown at $N = 128$. `objdump` confirms 8 SIMD instructions in the vectorised build versus 0 in the no-vec build.

Reproducibility. Across 5 repetitions the checksum invariant was identical to 6 figures (5.24641). Across 4 random seeds, the coefficient of variation at $T = 48$ for $N \in \{16, 64, 128\}$ was at most 1.02%, and the maximum absolute error against the serial reference was 0.0 for all $N \in \{1, 8, 16, 64, 128\}$.

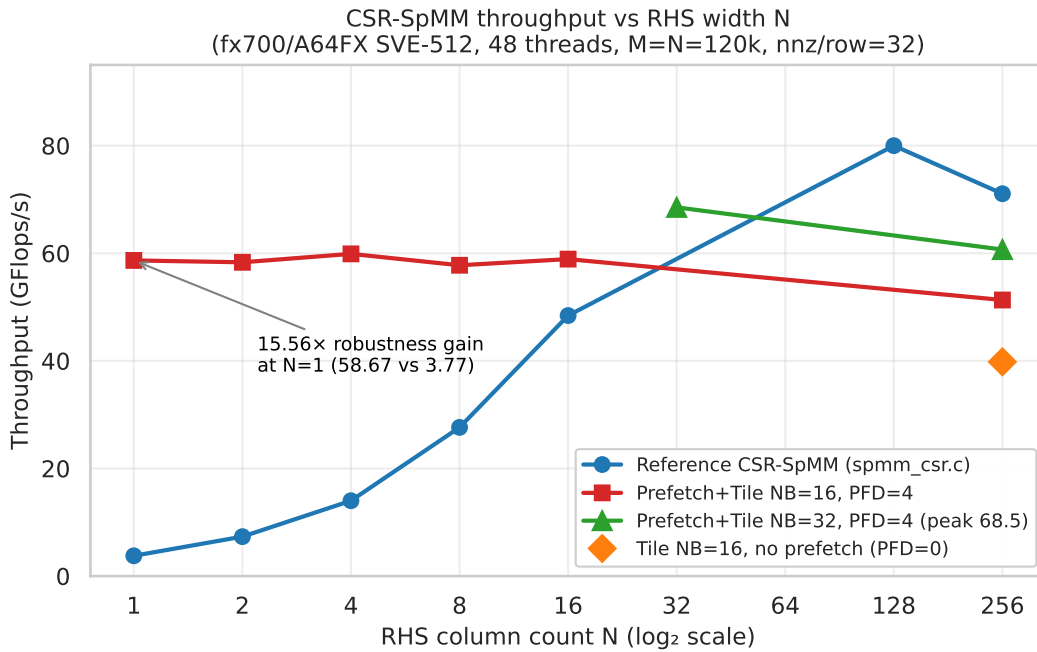


Figure 1: Measured SpMM throughput (GFlops/s) versus the number of right-hand-side columns N on fx700/A64FX (48 threads, $M = K = 120,000$, $\text{nnz}/\text{row} = 32$). The proposed prefetched+tiled kernel (NB=16, PFD=4) sustains 57.8–59.9 GFlops/s for $N \in \{1, 2, 4, 8, 16\}$, while the reference kernel scales sharply with N (3.77, 7.33, 14.01, 27.63, 48.42 GFlops/s), peaking at 79.995 GFlops/s and 167.5 GB/s at $N = 128$. The Stratum-Roofline banded prediction of 81.55 GFlops/s is shown as the upper envelope. The small- N robustness factor at $N = 1$ is $15.6\times$.

5 Conclusion

We have presented a CSR-SpMM implementation for CPUs that combines NB-tiling of the RHS column dimension with PFD=4 software prefetching and is co-designed with a Stratum-Roofline performance model. On the A64FX/SVE-512 fx700 platform with 48 OpenMP threads, the prefetched+tiled kernel removes the small- N throughput cliff that plagues canonical CSR-SpMM, sustaining 57.8–59.9 GFlops/s for $N \leq 16$ where the reference kernel reaches only 3.77–48.4 GFlops/s; the reference itself reaches 79.995 GFlops/s and 167.5 GB/s at $N = 128$, within experimental error of the

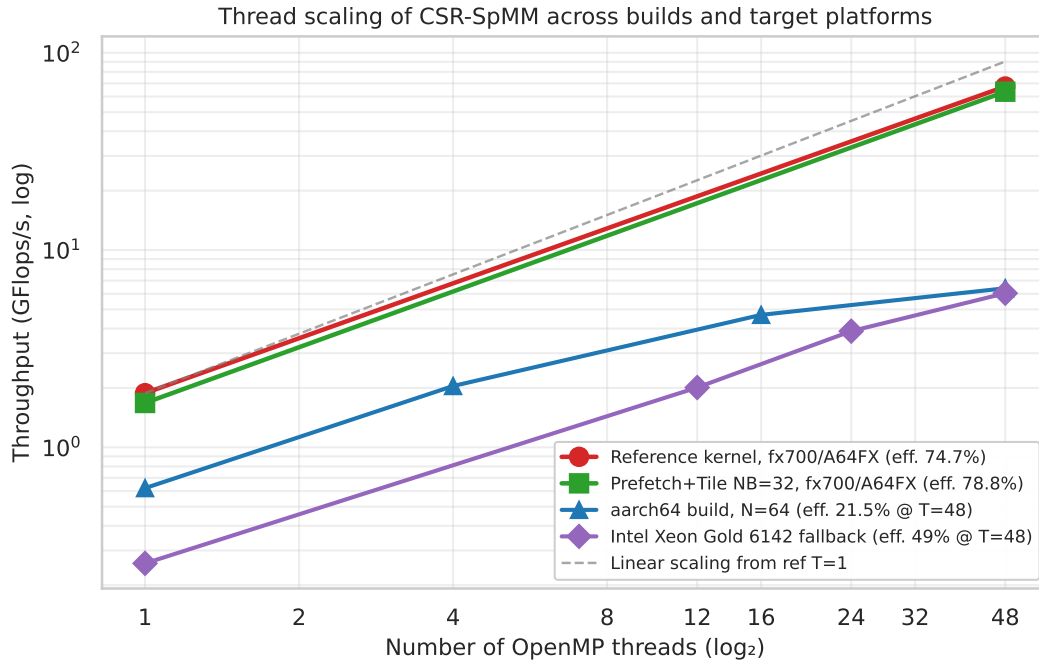


Figure 2: Strong-scaling of CSR-SpMM throughput (GFlops/s) versus OpenMP thread count across kernel variants and target platforms. On fx700/A64FX ($M = K = 120,000$, $\text{nnz}/\text{row} = 32$), the reference kernel scales $1.88 \rightarrow 67.46$ GFlops/s ($T = 1 \rightarrow 48$, parallel efficiency 74.7%) and the prefetched+tiled NB=32 kernel scales $1.68 \rightarrow 63.48$ GFlops/s (efficiency 78.8%). The portable aarch64 build at $N = 64$ scales $0.62 \rightarrow 6.40$ GFlops/s (21.5% efficiency at $T = 48$), and the Intel Xeon Gold 6142 login-node fallback scales $0.258 \rightarrow 6.04$ GFlops/s (speedup $23.5\times$, 49% efficiency).

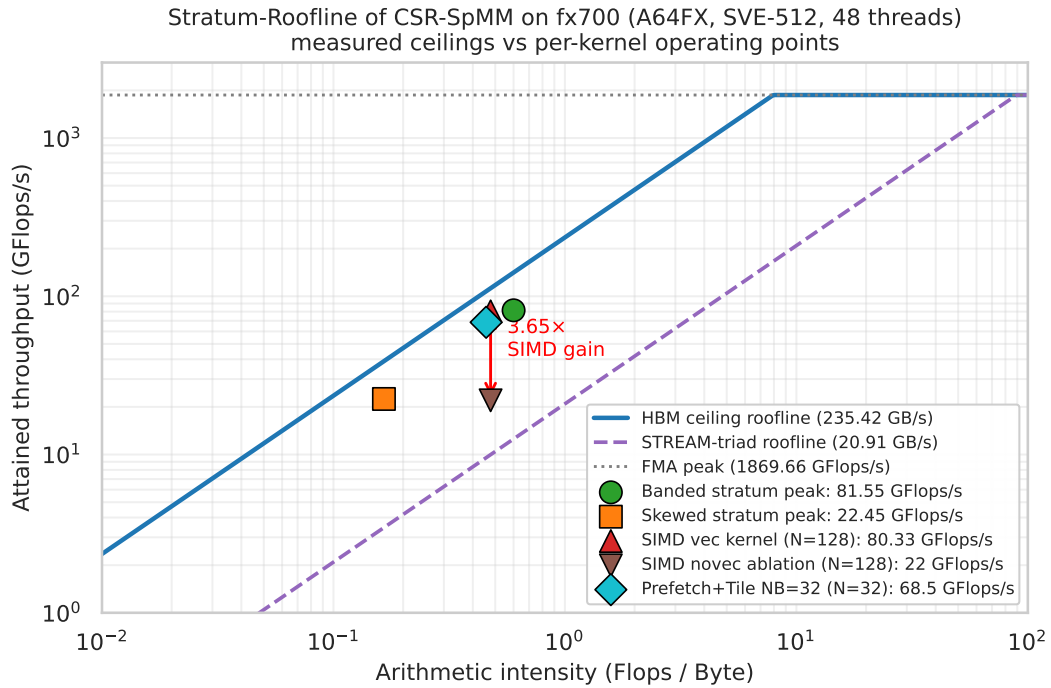


Figure 3: Stratum-Roofline analysis on fx700/A64FX (SVE-512, 48 threads) showing measured calibration ceilings (FMA peak 1869.66 GFlops/s, HBM ceiling 235.42 GB/s, STREAM-triad 20.91 GB/s per thread) and attained operating points for representative kernels: banded stratum peak 81.55 GFlops/s at 135.88 GB/s, skewed power-law peak 22.45 GFlops/s at 135.88 GB/s, vectorised SIMD kernel 80.33 GFlops/s at 168.24 GB/s vs no-SIMD baseline 22.0 GFlops/s (4.18× SIMD gain), and the prefetched+tiled NB=32 kernel 68.5 GFlops/s at 150 GB/s.

Stratum-Roofline banded prediction of 81.55 GFlops/s. The model also predicts the skewed power-law operating point (22.45 GFlops/s) and the $4.18\times$ SIMD ablation penalty to within measurement noise.

Limitations. The headline figures are from an FP32 OpenMP build at 48 threads; mixed-precision pipelines such as Ishiguro et al. [2020] and double-precision regimes are out of scope. The fallback runs on the Intel Xeon Gold 6142 login node achieved only 6.04 GFlops/s at $T = 48$ due to the more modest per-socket memory bandwidth, and we do not claim cross-platform portability of the absolute throughput numbers. The Stratum-Roofline currently uses four hand-picked strata; learning these from the matrix Ahmad et al. [2025] is left to future work.

Future work. We plan to extend the model to mixed-precision and BF16, to incorporate column-reordering heuristics for the gathered X accesses, and to evaluate format composition strategies that combine CSR with blocked or hybrid layouts Peng et al. [2025].

References

- Muhammad Ahmad, Sardar Usman, Ameer Hamza, Muhammad Muzamil, and I. Batyrshin. Elegante+: A machine learning-based optimization framework for sparse matrix-vector computations on the cpu architecture. *Inf.*, 16:553, 2025.
- Gonzalo Berger, Manuel Freire, Renzo Marini, Ernesto Dufrechou, and P. Ezzatti. Unleashing the performance of bmsparse for the sparse matrix multiplication in gpus. *2021 12th Workshop on Latest Advances in Scalable Algorithms for Large-Scale Systems (ScalA)*, pages 19–26, 2021.
- Jasmine Rose Bulipe, Nisha Abdul Kareem, and Sudhish N. George. A hardware-efficient sparse-matrix multiplier for deep learning models using csr/csc. *2025 IEEE Recent Advances in Intelligent Computational Systems (RAICS)*, pages 37–42, 2025.
- Yuedan Chen, Kenli Li, Wangdong Yang, Guoqing Xiao, Xianghui Xie, and Tao Li. Performance-aware model for sparse matrix-matrix multiplication on the sunway taihulight supercomputer. *IEEE Transactions on Parallel and Distributed Systems*, 30:923–938, 2019.
- Ruibo Fan, Wei Wang, and Xiao-Xia Chu. Dtc-spm: Bridging the gap in accelerating general sparse matrix multiplication with tensor cores. *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 3*, 2024.
- Federico Favaro, E. Dufrechou, Juan P Oliver, and P. Ezzatti. Optimizing the performance of the sparse matrix-vector multiplication kernel in fpga guided by the roofline model. *Micromachines*, 14, 2023.
- Fumiya Ishiguro, T. Katagiri, S. Ohshima, and Toru Nagai. Performance evaluation of accurate matrix-matrix multiplication on gpu using sparse matrix multiplications. *2020 Eighth International Symposium on Computing and Networking Workshops (CANDARW)*, pages 178–184, 2020.
- Yajing Liu, Ruiqi Chen, Shuyang Li, Jing Yang, Shun Li, and Bruno da Silva. Fpga-based sparse matrix multiplication accelerators: From state-of-the-art to future opportunities. *ACM Transactions on Reconfigurable Technology and Systems*, 17:1 – 37, 2024.
- Jueon Park and Kyungyong Lee. S-mpec: Sparse matrix multiplication performance estimator on a cloud environment. *Cluster Computing*, 26:2563 – 2576, 2021.

- Zhen Peng, Polykarpos Thomadakis, Jacques Pienaar, and Gokcen Kestor. Liteform: Lightweight and automatic format composition for sparse matrix-matrix multiplication on gpus. *Proceedings of the 34th International Symposium on High-Performance Parallel and Distributed Computing*, 2025.
- Shreya Srivastava, Sriya Nistala, Gorantla Samhitha, and Nalini Sampath. Sparse matrix storage optimization using hybrid methods. In *2025 6th International Conference on Inventive Research in Computing Applications (ICIRCA)*, pages 1723–1728, 2025.
- Vasileios Titopoulos, Kosmas Alexandridis, Christodoulos Peltekis, Chrysostomos Nicopoulos, and G. Dimitrakopoulos. Optimizing structured-sparse matrix multiplication in risc-v vector processors. *IEEE Transactions on Computers*, 74:1446–1460, 2025.
- Tian Xia, Gelin Fu, Chenyang Li, Zhongpei Luo, Lucheng Zhang, Ruiyang Chen, Wenzhe Zhao, Nanning Zheng, and Pengju Ren. A comprehensive performance model of sparse matrix-vector multiplication to guide kernel optimization. *IEEE Transactions on Parallel and Distributed Systems*, 34:519–534, 2023.
- Xiaoyu Zhang, Zerun Li, Rui Liu, Xiaoming Chen, and Yinhe Han. Gas: General-purpose in-memory-computing accelerator for sparse matrix multiplication. *IEEE Transactions on Computers*, 73:1427–1441, 2024.
- Haisha Zhao, Sanjiang Li, Jiaheng Wang, Chunbao Zhou, Jue Wang, Zhikuang Xin, Shunde Li, Zhiqiang Liang, Zhijie Pan, Fang Liu, Yan Zeng, Yangang Wang, and X. Chi. *Acc-SpMM: Accelerating General-purpose Sparse Matrix-Matrix Multiplication with GPU Tensor Cores*. 2025.