# 3D DRAM Based Application Specific Hardware Accelerator for SpMV

**Electrical & Computer ENGINEERING**

Fazle Sadi, Larry Pileggi, Franz Franchetti

**Contact: fsadi@cmu.edu**

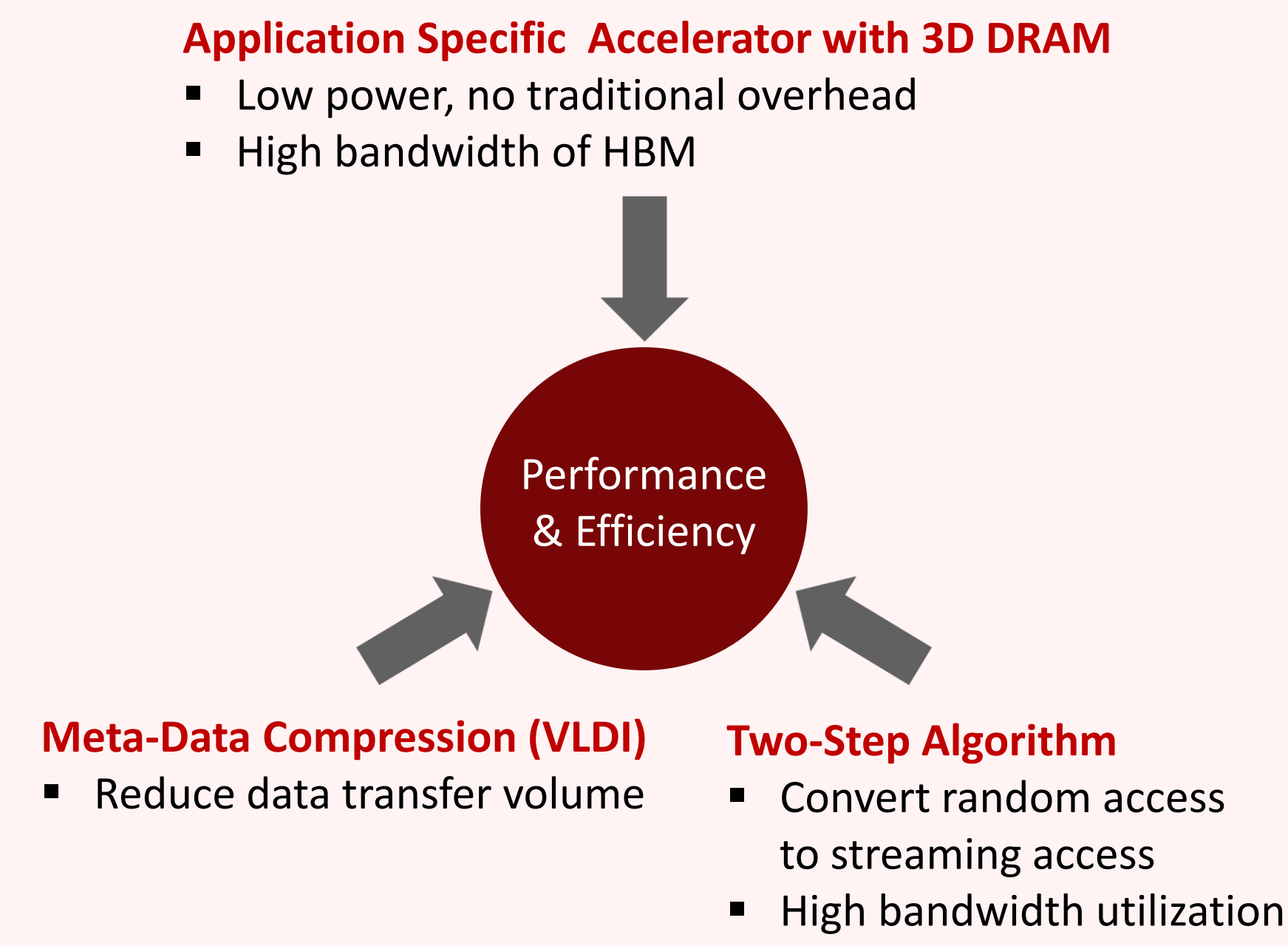**Carnegie Mellon University**

---

## Research Problem

Sparse Matrix-Vector multiplication (SpMV) is an important kernel for many applications. However, due to lack of data locality and low FLOP to memory access ratio, SpMV's performance and efficiency are very poor on regular architectures. Overcoming this problem warrants re-thinking of the way we do SpMV, both in terms of software & hardware.
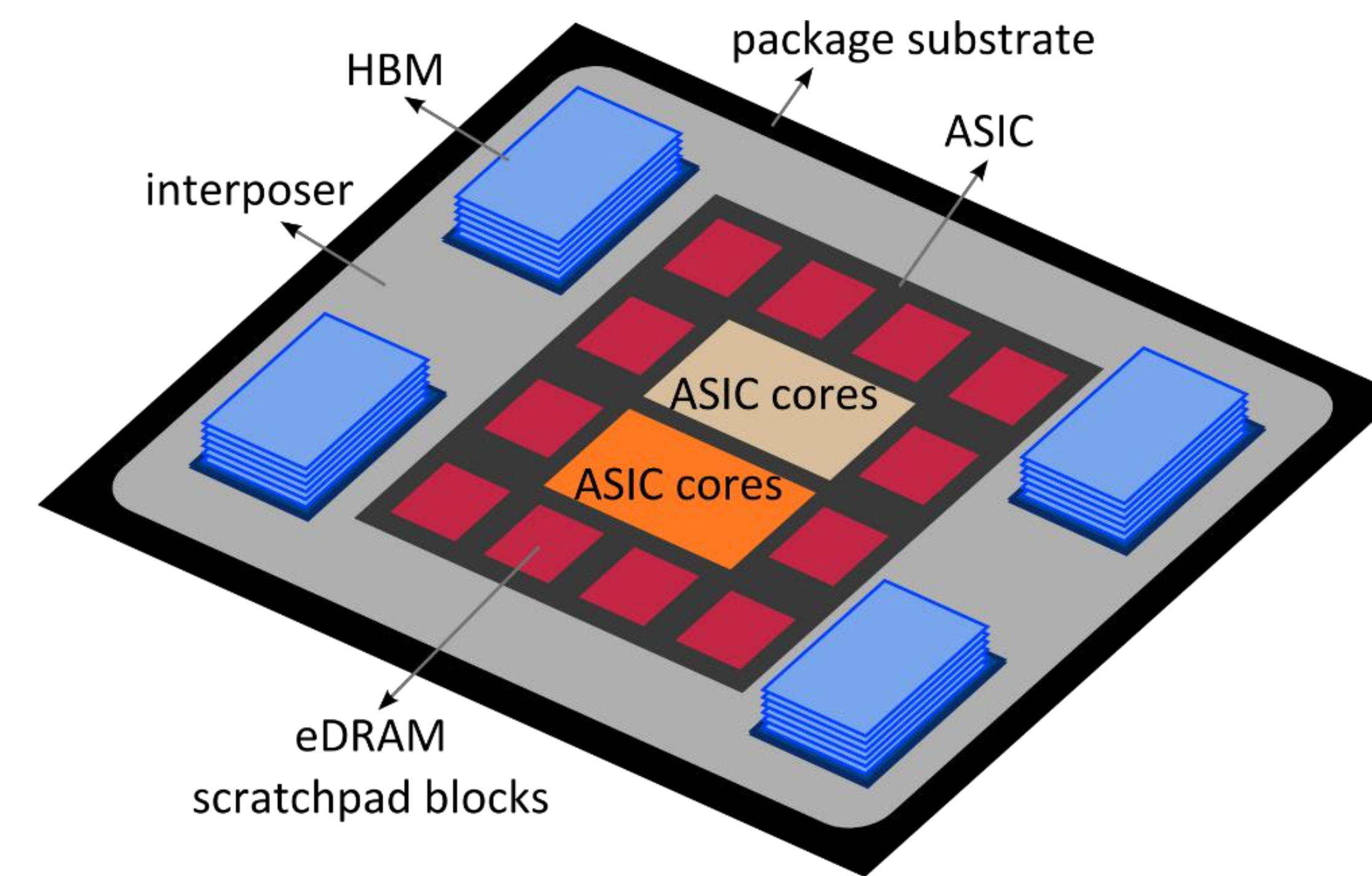


Peak CPU performance — Traditional memory hierarchy (CPU, Register, Cache, L1, L2, L3, Main Memory, RAM) — Sparse matrix — Huge address space

10% — Large SpMV performance — Expects temporal & spatial locality — No temporal or spatial locality

**SpMV does not benefit from traditional memory hierarchy.**

---

## Proposed Solution

We propose an algorithm-hardware co-optimization approach to gain significant performance and energy efficiency for SpMV.

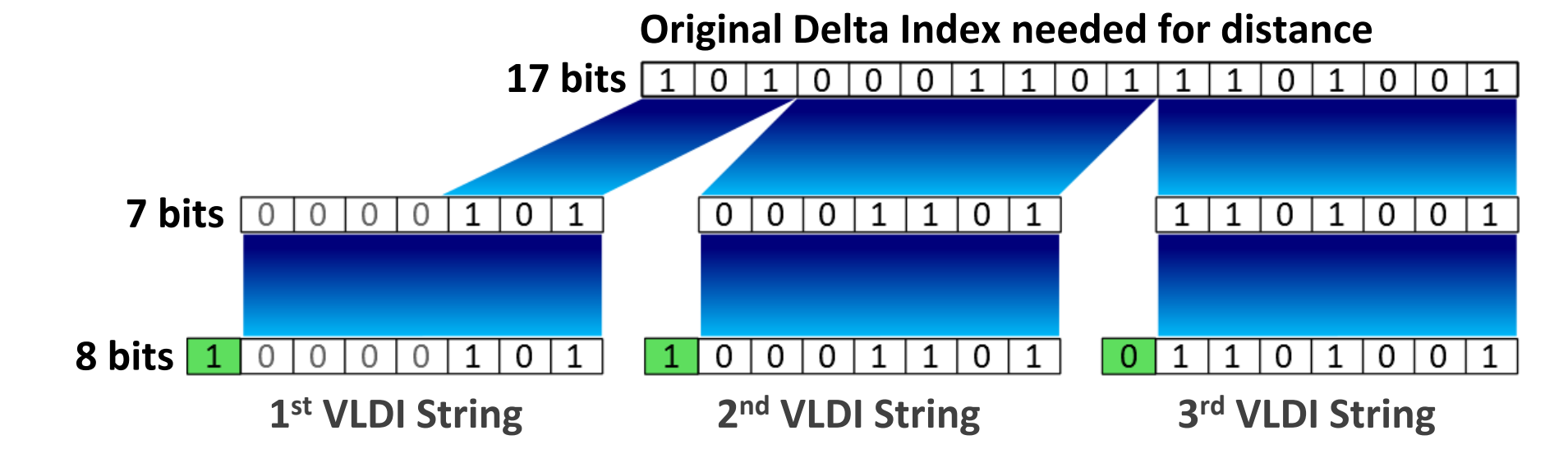**Application Specific Accelerator with 3D DRAM**
- Low power, no traditional overhead
- High bandwidth of HBM

Performance & Efficiency

**Meta-Data Compression (VLDI)**
- Reduce data transfer volume

**Two-Step Algorithm**
- Convert random access to streaming access
- High bandwidth utilization

---

**ALGORITHM AND HARDWARE CO-OPTIMIZED SPMV ACCELERATOR**



HBM — package substrate — ASIC — interposer — ASIC cores — ASIC cores — eDRAM scratchpad blocks
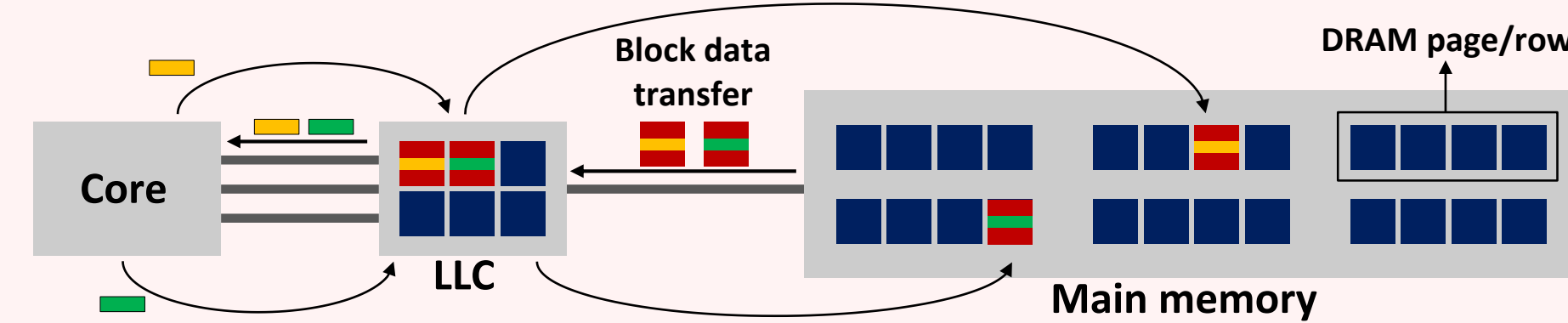
---

## Meta-Data Compression: VLDI

To reduce the meta-data storage and transfer cost, we use Variable Length Delta Index (VLDI) compression technique. Instead of absolute index, only the distances between non-zero elements are stored using VLDI strings.



Original Delta Index needed for distance

17 bits — 7 bits — 8 bits

1st VLDI String — 2nd VLDI String — 3rd VLDI String

**VLDI String Leading Bit**
- ☐ '1' means continuation
- ☐ '0' means termination

**It works because vector elements are**
- Generated/stored in order – *Step 1*
- Accessed sequentially – *Step 2*
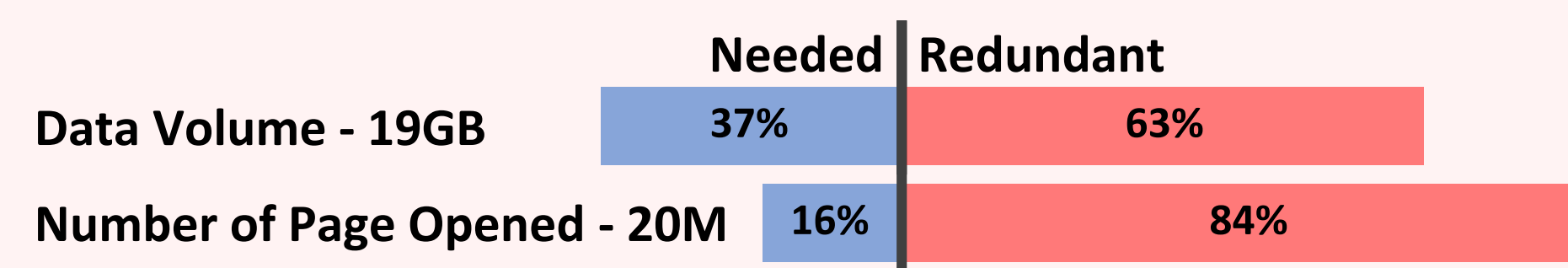
---

## Background

For every two computes on average SpMV requires a random access to dense vector which is generally too large for the last level cache (LLC). This results in random access to DRAM which affects the system performance and efficiency in various ways.



Core — LLC — Block data transfer — Main memory — DRAM page/row

1. Low bandwidth utilization ⟹ Poor performance
2. Excess volume & page opening ⟹ Poor efficiency

### An Example

80M x 80M matrix, 3 NNZ per row, 1KB DRAM page, 64B cache block, double precision data

| | Needed | Redundant |
|---|---|---|
| Data Volume - 19GB | 37% | 63% |
| Number of Page Opened - 20M | 16% | 84% |

**Solution: Convert random accesses to streaming access.**

---

## Will HBM with COTS work?

Probably not.

The key problem with SpMV on commercial off-the-shelf (COTS) architectures is the traditional memory hierarchy. Adding 3D stacked DRAM, such as High Bandwidth Memory (HBM), will only provide more DRAM bandwidth. Current SpMV implementations on COTS platforms will still have the same poor bandwidth utilization and other pertinent issues. Therefore, any significant improvement of performance or efficiency is unlikely as long as SpMV is implemented on any architecture which is built on traditional memory hierarchy.

---

## Two-Step Algorithm



Source Matrix — Source Vector — X — Multiplication & Partial Addition — **Step 1** — Intermediate Vectors (Sparse) — Big Multi-Way Merge — **Step 2** — Target Vector
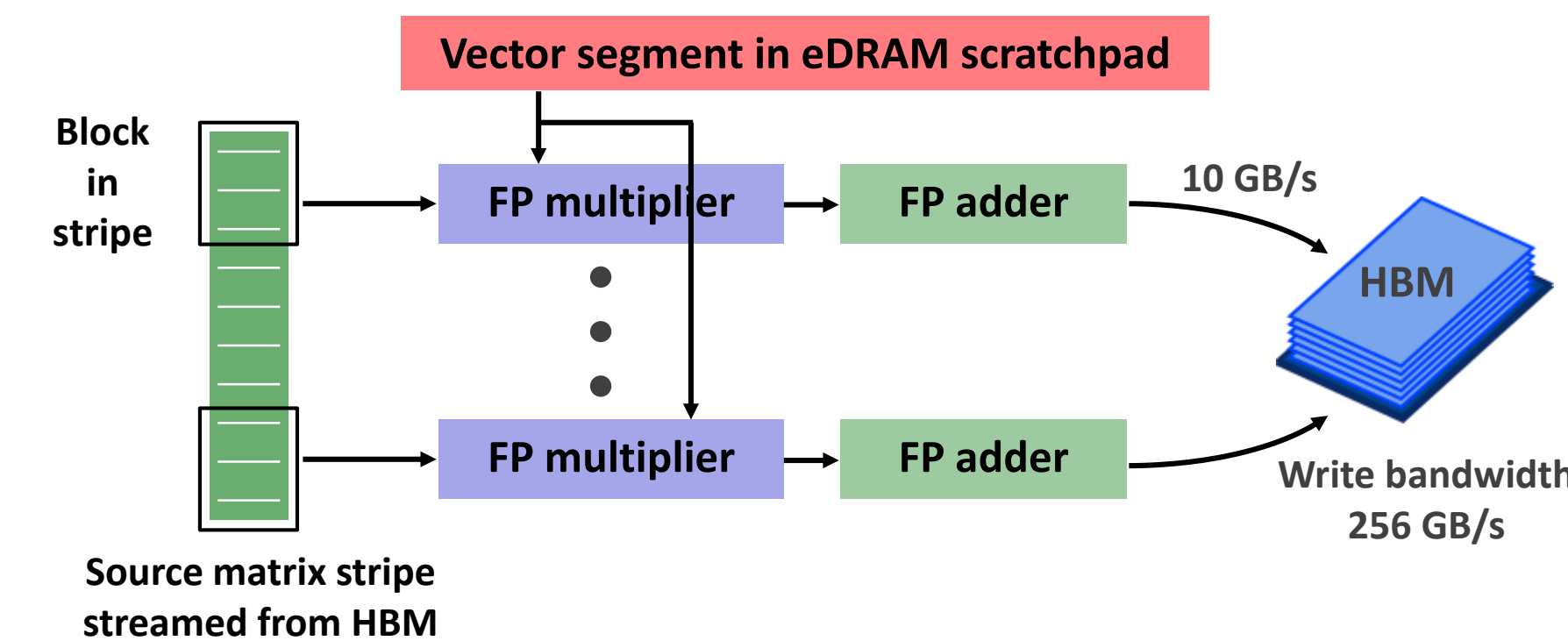
Each stripe in CSR format — Data value + Index (Meta-data)

- ☐ Conversion of all DRAM random accesses into sequential accesses.
- ☐ Dedicated scalable multi-way merge network is required.
- ☐ ASIC merge is capable of merging large number of lists (e.g. 1k) with high throughput (1.5B elements/s)
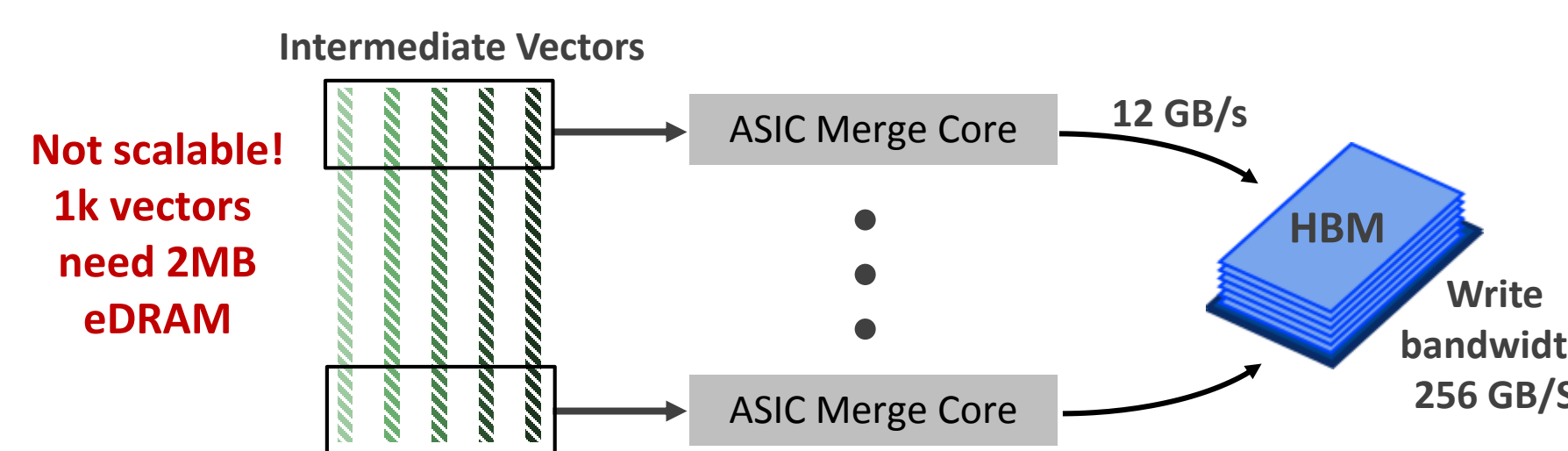
### Step 1: Implementation and Parallelization

We can easily scale the process in step 1 to utilize full HBM bandwidth by blocking the matrix stripes.



Vector segment in eDRAM scratchpad — Block in stripe — FP multiplier — FP adder — 10 GB/s — HBM — Source matrix stripe streamed from HBM — Write bandwidth 256 GB/s
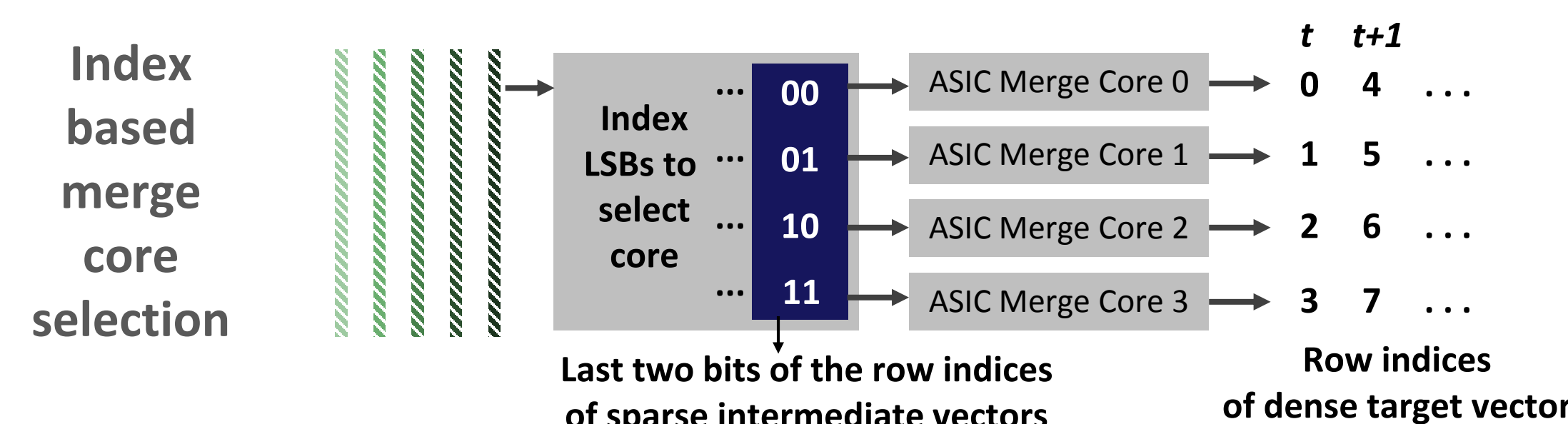
### Step 2: Implementation and Parallelization

It is not possible to scale the global multi-way merge in step 2 just by blocking the intermediate sparse vectors.



Intermediate Vectors — **Not scalable! 1k vectors need 2MB eDRAM** — ASIC Merge Core — 12 GB/s — HBM — ASIC Merge Core — Write bandwidth 256 GB/S
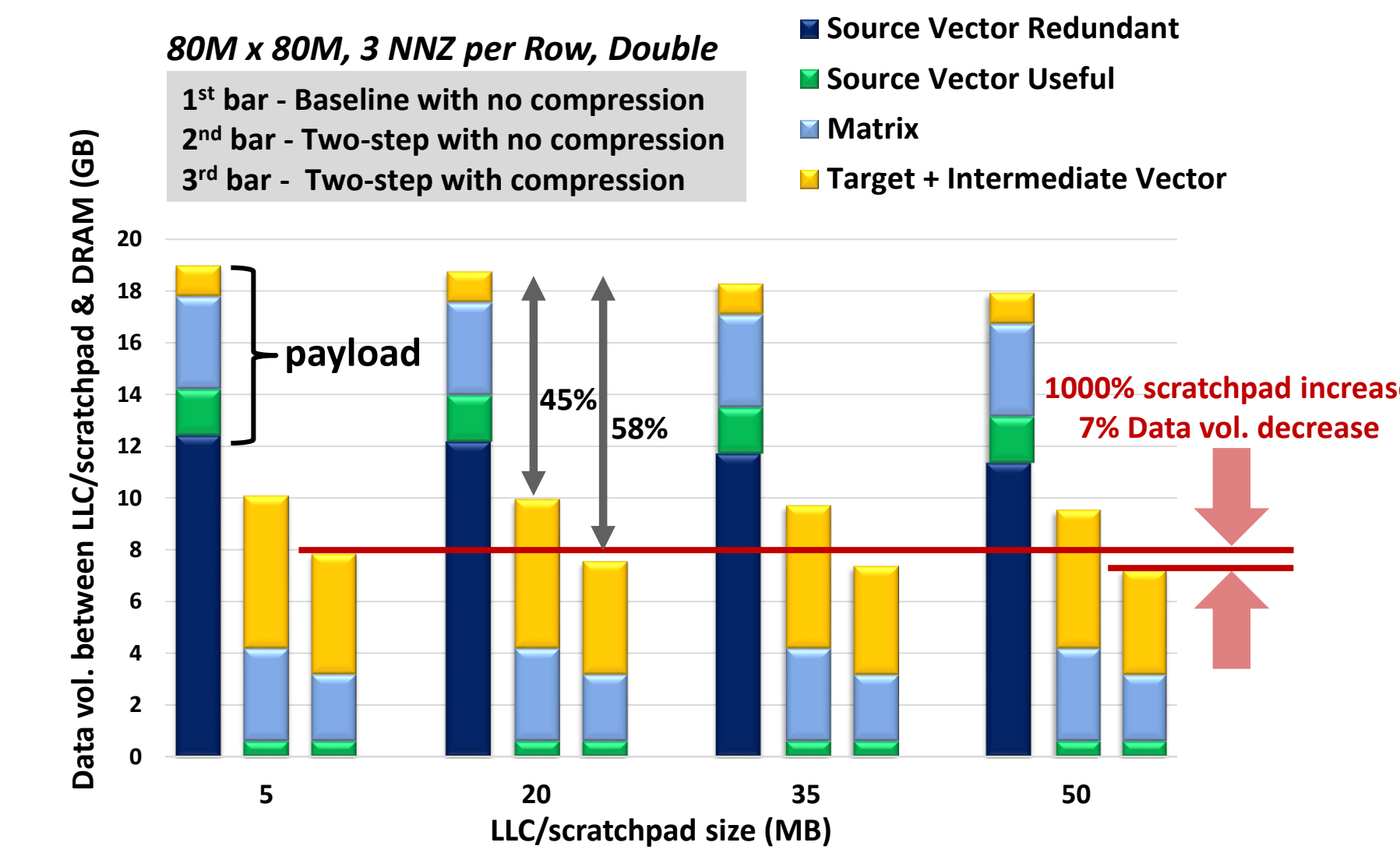
To make merge scalable, we assign the merging task to different merge cores depending on few LSBs of the indices. This novel technique works for SpMV as it is guaranteed that the final target dense vector will have data for each index.
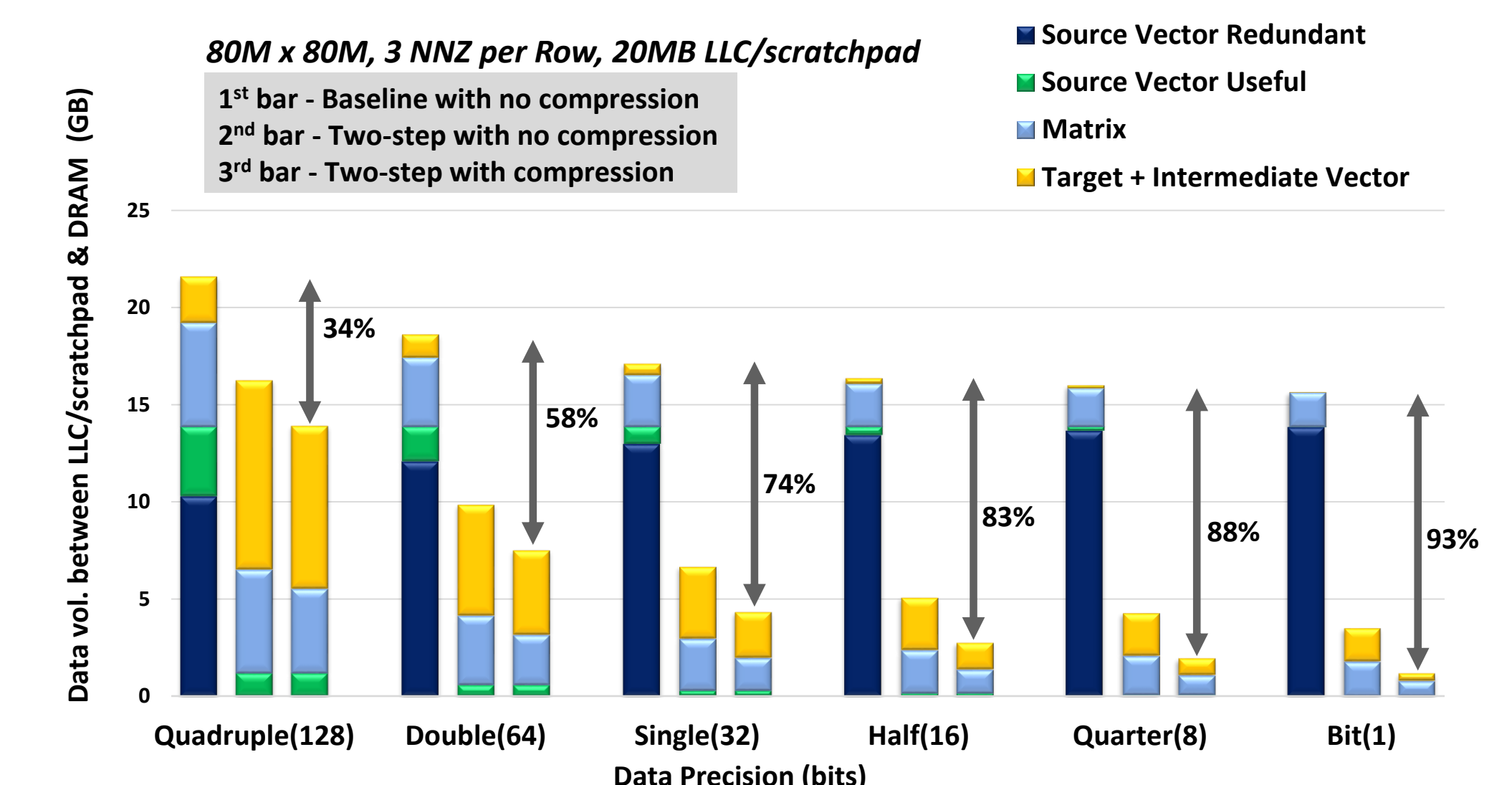


Index based merge core selection — Index LSBs to select core — 00 → ASIC Merge Core 0 → t: 0, t+1: 4 ... — 01 → ASIC Merge Core 1 → 1, 5 ... — 10 → ASIC Merge Core 2 → 2, 6 ... — 11 → ASIC Merge Core 3 → 3, 7 ...

Last two bits of the row indices of sparse intermediate vectors — Row indices of dense target vector

---

## Experimental Results

### Data Volume vs LLC



80M x 80M, 3 NNZ per Row, Double

1st bar - Baseline with no compression
2nd bar - Two-step with no compression
3rd bar - Two-step with compression

Legend: Source Vector Redundant, Source Vector Useful, Matrix, Target + Intermediate Vector

Data vol. between LLC/scratchpad & DRAM (GB) vs LLC/scratchpad size (MB): 5, 20, 35, 50

payload — 45% — 58% — 1000% scratchpad increase 7% Data vol. decrease

### Data Volume vs Precision



80M x 80M, 3 NNZ per Row, 20MB LLC/scratchpad

1st bar - Baseline with no compression
2nd bar - Two-step with no compression
3rd bar - Two-step with compression

Legend: Source Vector Redundant, Source Vector Useful, Matrix, Target + Intermediate Vector

Data vol. between LLC/scratchpad & DRAM (GB) vs Data Precision (bits): Quadruple(128) 34%, Double(64) 58%, Single(32) 74%, Half(16) 83%, Quarter(8) 88%, Bit(1) 93%
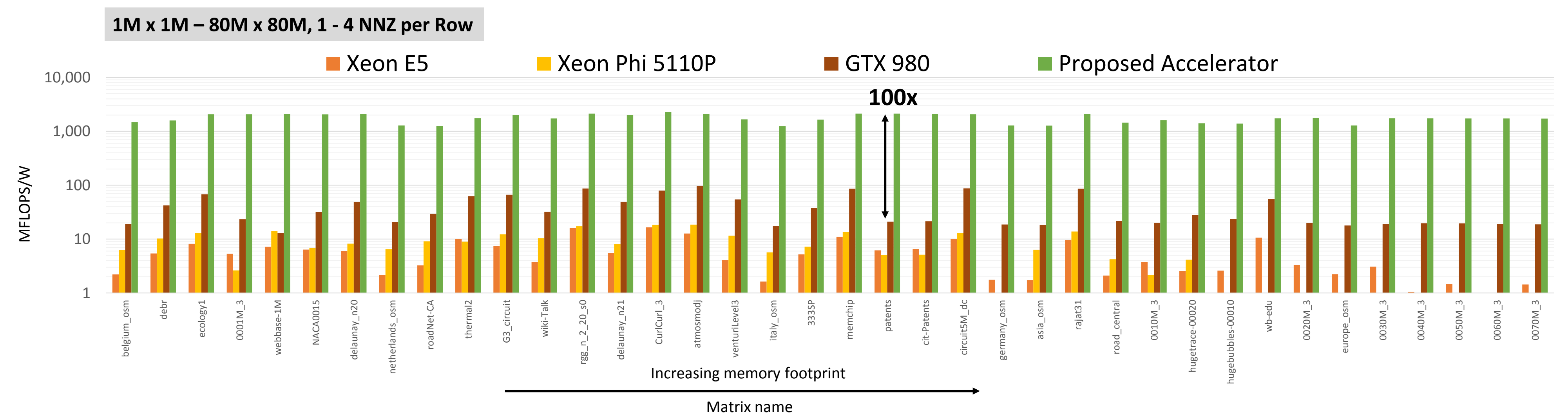
- ☐ Even though we increase the payload with Two-step, the overall data transfer volume is significantly reduced.
- ☐ Large eDRAM size does not significantly reduce data transfer volume. With efficient merge network small (i.e. low cost) scratchpad will suffice.
- ☐ For lower data precisions, the compression technique becomes more effective as meta-data becomes more dominant.

### Energy Efficiency Comparison

- ☐ CPU(Xeon E5) + MKL : 22nm, 30MB LLC, 102.4GB/s Peak BW
- ☐ Co-processor (Xeon Phi) + MKL: 22nm, 30MB LLC, 352GB/s Peak BW
- ☐ GPU (GTX 980) + cuSPARSE: 28nm, 2MB LLC, GDDR5 224GB/s Peak BW
- ☐ Proposed Accelerator + Two-Step: 28nm, 20MB scratchpad, HBM 512GB/s Peak BW



1M x 1M – 80M x 80M, 1 - 4 NNZ per Row

Legend: Xeon E5, Xeon Phi 5110P, GTX 980, Proposed Accelerator

MFLOPS/W vs Matrix name — 100x — Increasing memory footprint

Large sparse matrices with high sparsity and various sizes from University of Florida's collection have been used to compare the energy efficiency of the proposed system against state of the art architectures. It can be seen that the proposed accelerator can achieve up to 100x more energy efficiency than the best performing COTS architecture.

---

## Acknowledgements