+

# Lecture 15: Sparse Matrix-Vector Multiplication (SpMV)

## Helen Xu
## hxu615@gatech.edu

Georgia Tech College of Computing
**School of Computational Science and Engineering**

# Announcements

- HW4 deadline postponed to next Tu, March 5 @ 5pm

- Please sign up for project groups by next Monday, March 4, by 5pm

- Project proposal writeup / slides due Tu, March 12 (2 weeks) - we may need to schedule some presentations outside of class time - we will release an initial schedule after groups are formed.

- HW5 out - deadline postponed to Tu, Mar 26 (after the spring break)

- As a reminder, please don't use AI tools (e.g., ChatGPT, etc.). Also, you are welcome to discuss with other students about the HW, but please don't look directly at or copy other's code.

# 1990 Nobel Prize in Economics



Our models strained the computer capabilities of the day [1950s]. I observed that **most of the coefficients in our matrices were zero**; i.e. , the nonzeros were "sparse" in the matrix
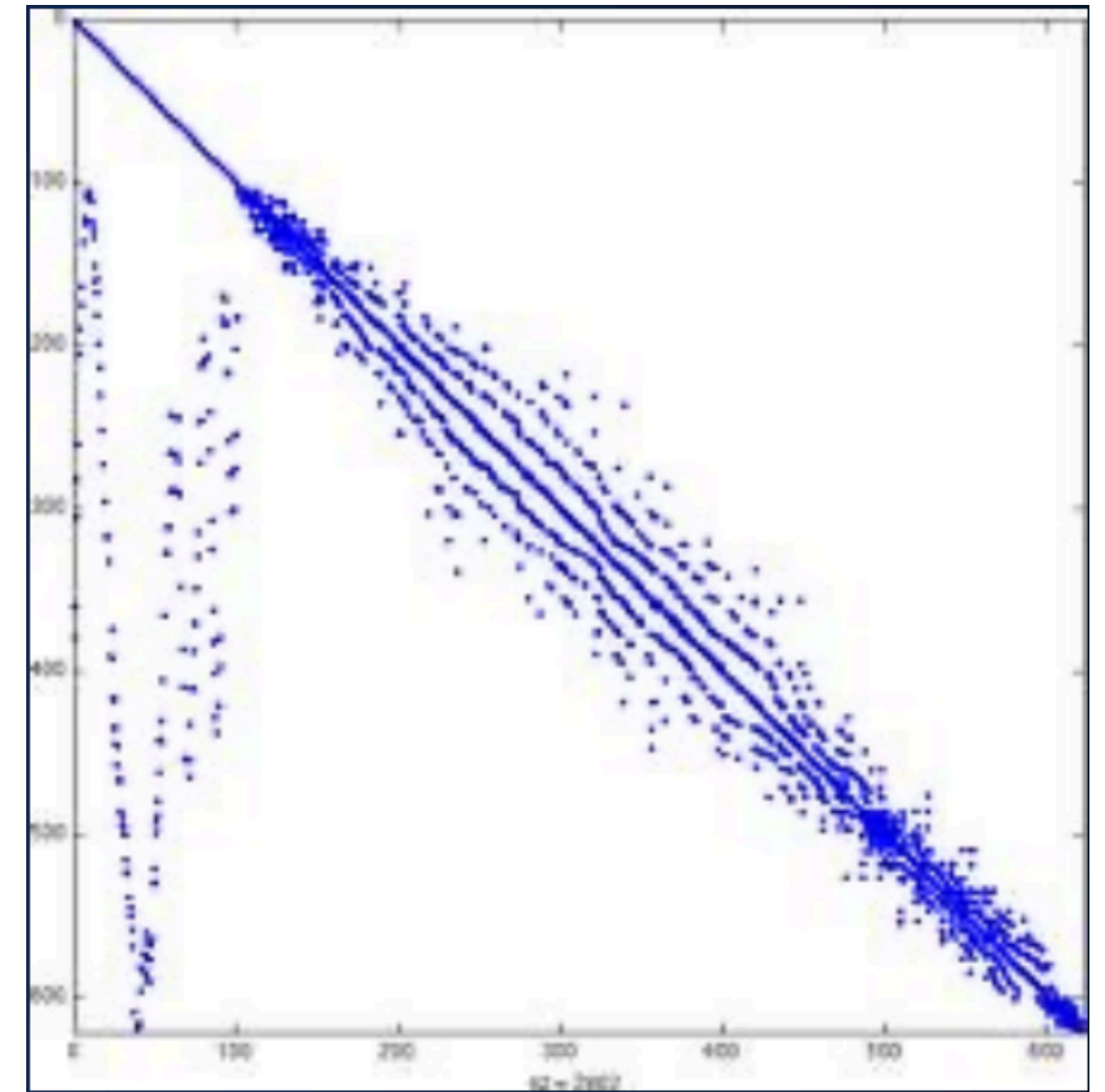
- Harry Markowitz

# What is a sparse matrix?

A matrix with **primarily zeros** (>> 90%).

Representing them using dense data structures **wastes memory and computation**.

Sparse matrices arise in many applications
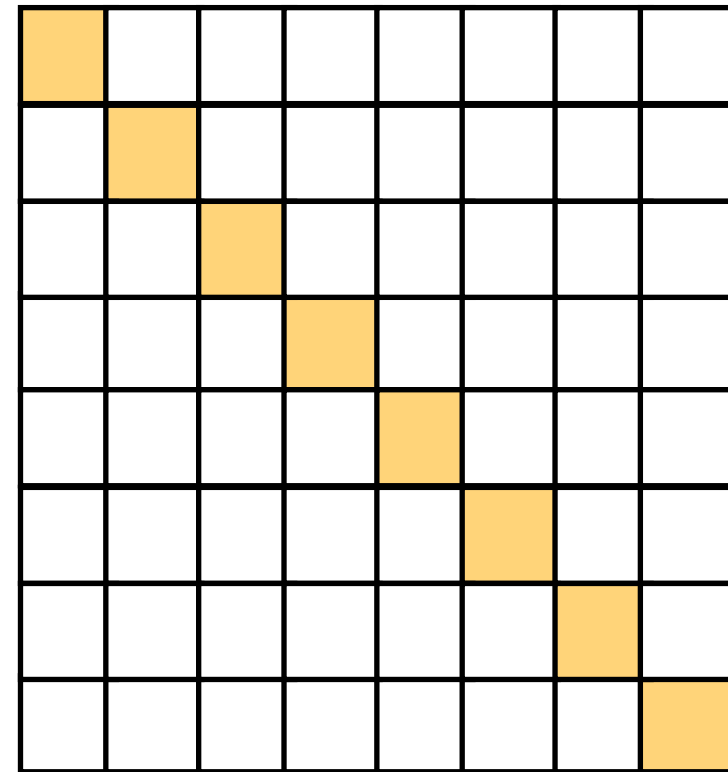- Simulating climate
- Analyzing images (photos, MRIs,…)
- Web page ranking for search
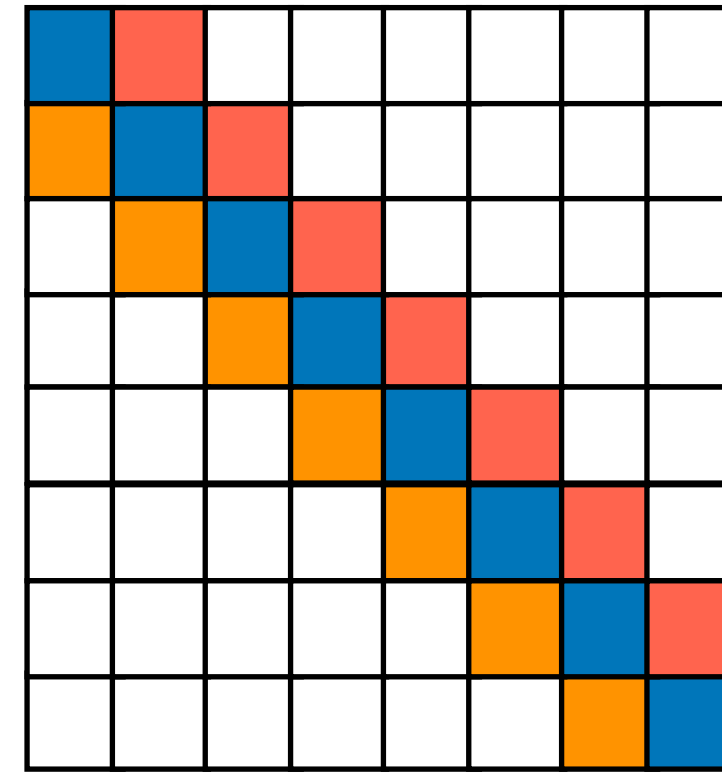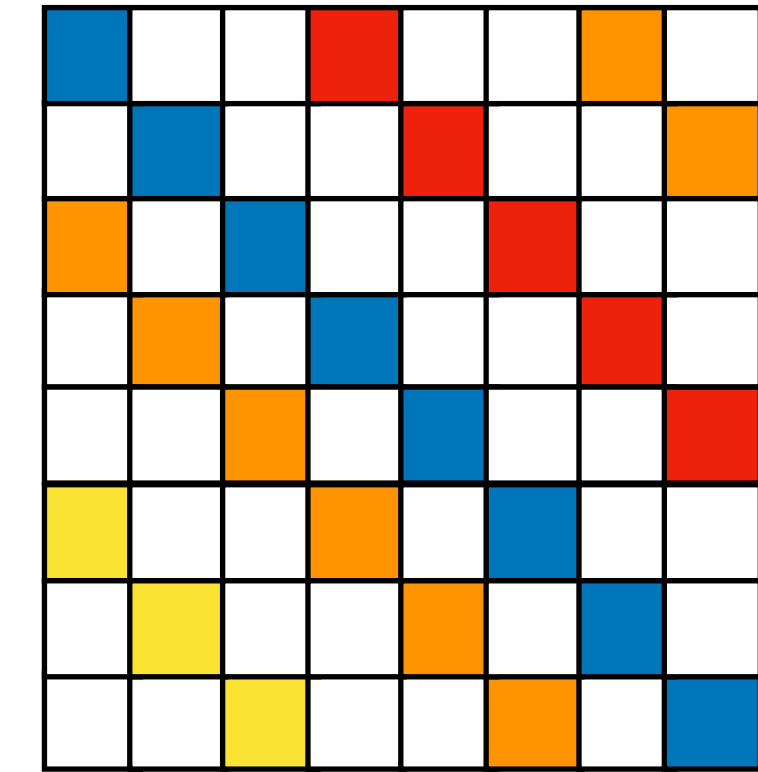- Graphs, including Graph Neural Nets



**https://sparse.tamu.edu/**
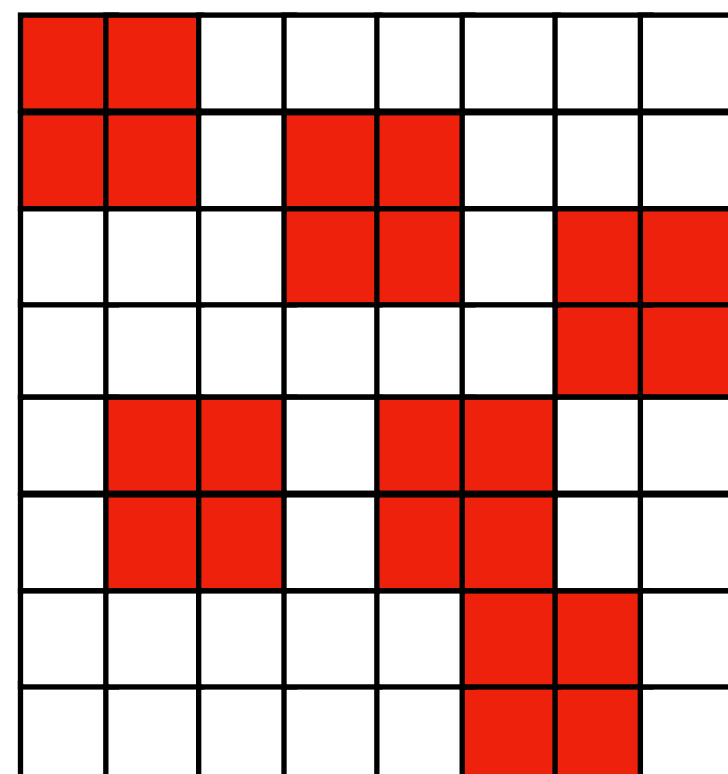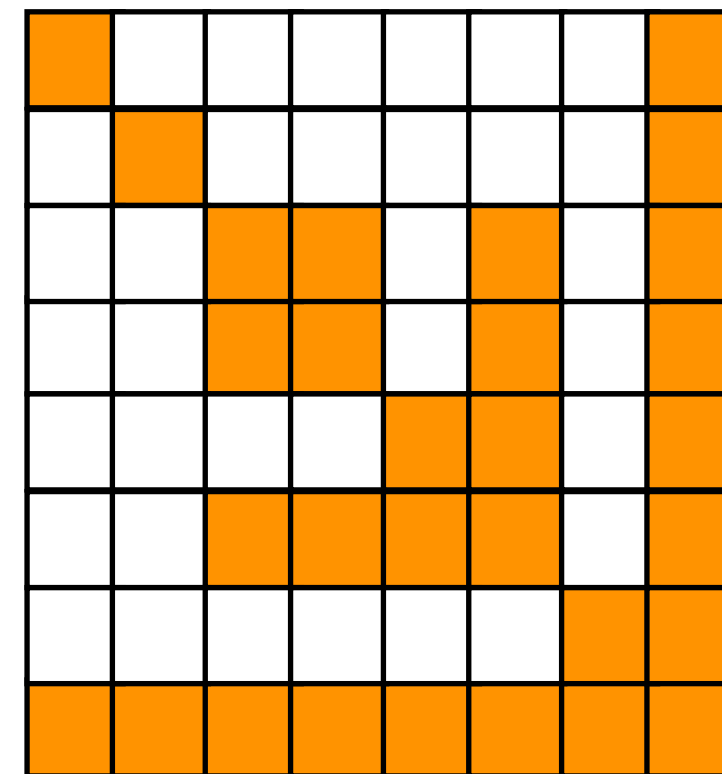
# Examples of sparse matrices
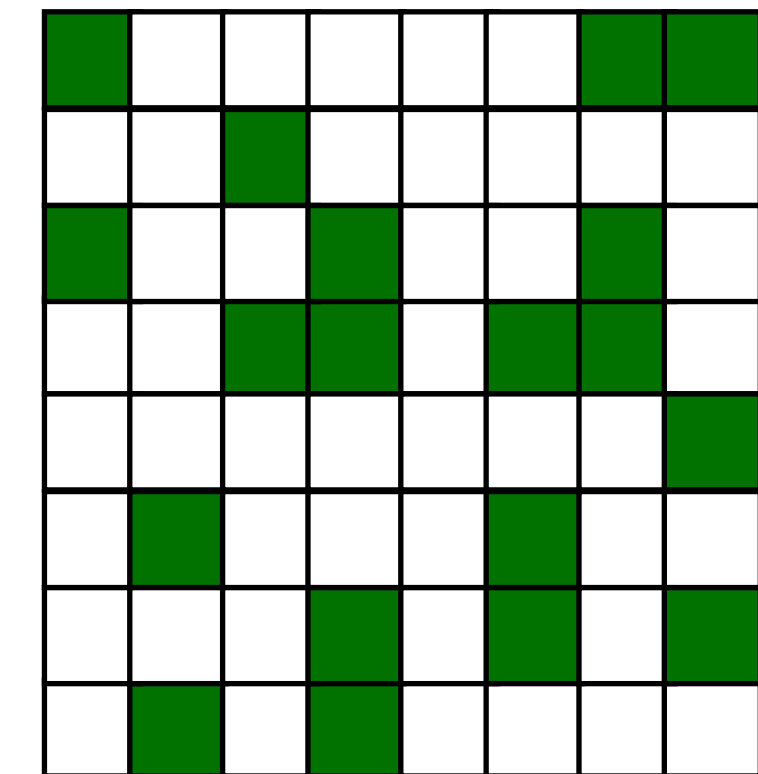
**Diagonal**

**Tridiagonal**

**"Generalized diagonal"**

**Block matrix
(2x2 dense blocks)**

**Symmetric**

**Irregular**

# Sparse matrices are everywhere

**Internet connectivity**  **Structural design**  **Linear programming**

# Recommendation matrix

**Products**

**Population**



**Ratings**

# Image Segmentation

Image segmentation – Identify the **object boundaries** in an image

Compute **affinity matrix** – Is it likely that two pixels belong to the same object?



Figure 7. Selected image contours

"Efficient, High-Quality Image Contour Detection" Catanzaro, Su, Sundaram, Lee, Murphy, Keutzer, International Conference on Computer Vision, September 2009

# More applications

Graphs
- Google's PageRank (originally an eigen-problem on the web adjacency matrix)
- Transportation network analysis

Text analysis
- Latent Semantic Indexing finds topics in a document corpus by doing a singular value decomposition of a bag-of-words matrix

Scientific and engineering
- Solving differential equations (climate modeling, etc.)
- Optimization problems

And many more…

# Sparse Matrix Formats

# Adjacency matrix representation

Any sparse matrix representation can be used for sparse graphs, and vice versa.

# Coordinate representation (COO)



(0, 0, 12)
(0, 2, 26)
(2, 1, 19)
(3, 1, 14)
(3, 3, 7)

row + column index +
weight per nonzero
(easy to build / modify)

# Adjacency lists



Source stored implicitly

# Compressed Sparse Row (CSR)

Compressed sparse row (CSR) = **cache-efficient adjacency lists**



|   | 0 | 1 | 2 | 3 | 4 |
|---|---|---|---|---|---|
| **ptr** | 0 | 2 | 2 | 3 | 5 | **(row starts in CSR)**

Index into adjacency array

| **ind** | 0 | 2 | 1 | 1 | 3 | **(column ids in CSR)**

Adjacencies

| **val** | 12 | 26 | 19 | 14 | 7 | **(numerical values in CSR)**

Weights

# Directed vs Undirected

An undirected graph corresponds to a symmetric matrix - only need to store **half**



|   | **0** | **1** | **2** | **3** |
|---|---|---|---|---|
| **0** | 12 |   | 26 |   |
| **1** |   |   | 19 | 14 |
| **2** | 26 | 19 |   |   |
| **3** |   | 14 |   | 7 |

|   | **0** | **1** | **2** | **3** | **4** | |
|---|---|---|---|---|---|---|
| **ptr** | 0 | 2 | 2 | 3 | 5 | **(row starts in CSR)** |
| **ind** | 0 | 2 | 1 | 1 | 3 | **(column ids in CSR)** |
| **val** | 12 | 26 | 19 | 14 | 7 | **(numerical values in CSR)** |

Index into adjacency array

Adjacencies

Weights

# Compressed Sparse Row (CSR) Storage



CSR has:
- Size **nnz** = number of nonzeros
- Array of the nonzero **values** (val) of size nnz
- Array of the column **indices** (ind) for each value of size nnz
- Array of row start **pointers** (ptr) of size n = number of rows

# Other Storage Formats

**Compressed Sparse Row (CSR) is the most common** and our focus today

Others include
- Compressed Sparse Column (CSC)
- Diagonal (DIAG): store main diagonal as 1D array; or diagonal bands as 2D (padded)
- Symmetric: store only ½ the array (indexing more complicated)
- **Blocked**: store each block contiguously
  - **Register blocked**: blocks are small and dense, avoid indexes within blocks
  - **Cache blocked**: blocks are large and themselves sparse

…and many other specialized formats!

# Serial SpMV

# Serial SpMV in CSR

**SpMV: Sparse Matrix-Vector Multiplication**

**Representation of A**



Compute **y(i) = y(i) + A(i, j) * x(j)**

In CSR

# Serial SpMV in CSR

**SpMV: Sparse Matrix-Vector Multiplication**

**Representation of A**



In CSR

Compute **y(i) = y(i) + A(i, j) * x(j)**

```
for each row i:
  for k = ptr[i] to ptr[i+1] do
    y[i] += val[k] * x[ind[k]]
```

- No reuse in A
- Maximum reuse in y as written
- Reuse in x?

# Parallel SpMV

# SpMV in CSR: OpenMP Parallel

```
#pragma omp parallel num_threads(thread_num)
{
#pragma omp for private(j, i, tmp) schedule(static)
  for (int i = 0; i < m; i++) {
    for (int j = ptr[i]; j < ptr[i+1]; j++) {
      tmp = ind[j];
      y[i] += val[j] * x[tmp]
    }
  }
}
```

# SpMV in CSR: OpenMP Parallel

```
#pragma omp parallel num_threads(thread_num)
{
#pragma omp for private(j, i, tmp) schedule(dynamic)
  for (int i = 0; i < m; i++) {
    for (int j = ptr[i]; j < ptr[i+1]; j++) {
      tmp = ind[j];
      y[i] += val[j] * x[tmp]
    }
  }
}
```

**May prefer dynamic because rows are not uniformly sized**



x

y    A

# SpMV for 8-wide SIMD

```c
void avx2_csr_spmv( float *A, int32_t *nIdx, int32_t **indices, float *x, int32_t m, float *y) {
  int32_t A_offset = 0;
  for(int32_t i = 0; i < m; i++) {
    int32_t nElem = nIdx[i]; float t = 0.0f;
    __m256 vT = _mm256_set_ps(0.0f,0.0f,0.0f,0.0f,0.0f,0.0f,0.0f,0.0f);
    int32_t smLen = nElem - (nElem & 7);
    for(int32_t j = 0; j < smLen; j+=8) {
      __m256i vIdx = _mm256_loadu_si256((__m256i*)&(indices[i][j]));
      __m256 vX = _mm256_i32gather_ps((float const*)x,vIdx,4);
      __m256 vA = _mm256_loadu_ps(&A[A_offset + j]);
      vT = _mm256_add_ps(vT, _mm256_mul_ps(vX,vA));
    }
    t += sum8(vT);
    for(int32_t j = smLen; j < nElem; j++) {
      int32_t idx = indices[i][j];
      t += x[idx]*A[A_offset + j];
    }
    y[i] = t;
    A_offset += nElem;
  }
}
```

# SpMV for 8-wide SIMD

```
void avx2_csr_spmv( float *A, int32_t *nIdx, int32_t **indices, float *x, int32_t m, float *y) {
  int32_t A_offset = 0;
  for(int32_t i = 0; i < m; i++) {
    int32_t nElem = nIdx[i]; float t = 0.0f;
    __m256 vT = _mm256_set_ps(0.0f,0.0f,0.0f,0.0f,0.0f,0.0f,0.0f,0.0f);
    int32_t smLen = nElem - (nElem & 7);
    for(int32_t j = 0; j < smLen; j+=8) {
      __m256i vIdx = _mm256_loadu_si256((__m256i*)&(indices[i][j]));
      __m256 vX = _mm256_i32gather_ps((float const*)x,vIdx,4);
      __m256 vA = _mm256_loadu_ps(&A[A_offset + j]);
      vT = _mm256_add_ps(vT, _mm256_mul_ps(vX,vA));
    }
    t += sum8(vT);
    for(int32_t j = smLen; j < nElem; j++) {
      int32_t idx = indices[i][j];
      t += x[idx]*A[A_offset + j];
    }
    y[i] = t;
    A_offset += nElem;
  }
}
```

Iterate over rows

# SpMV for 8-wide SIMD

```
void avx2_csr_spmv( float *A, int32_t *nIdx, int32_t **indices, float *x, int32_t m, float *y) {
  int32_t A_offset = 0;
  for(int32_t i = 0; i < m; i++) {
    int32_t nElem = nIdx[i]; float t = 0.0f;
    __m256 vT = _mm256_set_ps(0.0f,0.0f,0.0f,0.0f,0.0f,0.0f,0.0f,0.0f);
    int32_t smLen = nElem - (nElem & 7);
    for(int32_t j = 0; j < smLen; j+=8) {
      __m256i vIdx = _mm256_loadu_si256((__m256i*)&(indices[i][j]));
      __m256 vX = _mm256_i32gather_ps((float const*)x,vIdx,4);
      __m256 vA = _mm256_loadu_ps(&A[A_offset + j]);
      vT = _mm256_add_ps(vT, _mm256_mul_ps(vX,vA));
    }
    t += sum8(vT);
    for(int32_t j = smLen; j < nElem; j++) {
      int32_t idx = indices[i][j];
      t += x[idx]*A[A_offset + j];
    }
    y[i] = t;
    A_offset += nElem;
  }
}
```

Iterate over rows

For each 8 nonzeroes

# SpMV for 8-wide SIMD

```c
void avx2_csr_spmv( float *A, int32_t *nIdx, int32_t **indices, float *x, int32_t m, float *y) {
  int32_t A_offset = 0;
  for(int32_t i = 0; i < m; i++) {
    int32_t nElem = nIdx[i]; float t = 0.0f;
    __m256 vT = _mm256_set_ps(0.0f,0.0f,0.0f,0.0f,0.0f,0.0f,0.0f,0.0f);
    int32_t smLen = nElem - (nElem & 7);
    for(int32_t j = 0; j < smLen; j+=8) {
      __m256i vIdx = _mm256_loadu_si256((__m256i*)&(indices[i][j]));
      __m256 vX = _mm256_i32gather_ps((float const*)x,vIdx,4);
      __m256 vA = _mm256_loadu_ps(&A[A_offset + j]);
      vT = _mm256_add_ps(vT, _mm256_mul_ps(vX,vA));
    }
    t += sum8(vT);
    for(int32_t j = smLen; j < nElem; j++) {
      int32_t idx = indices[i][j];
      t += x[idx]*A[A_offset + j];
    }
    y[i] = t;
    A_offset += nElem;
  }
}
```

Iterate over rows

For each 8 nonzeroes

Gather X, load A, multiply, add to y

# SpMV for 8-wide SIMD

```c
void avx2_csr_spmv( float *A, int32_t *nIdx, int32_t **indices, float *x, int32_t m, float *y) {
  int32_t A_offset = 0;
  for(int32_t i = 0; i < m; i++) {
    int32_t nElem = nIdx[i]; float t = 0.0f;
    __m256 vT = _mm256_set_ps(0.0f,0.0f,0.0f,0.0f,0.0f,0.0f,0.0f,0.0f);
    int32_t smLen = nElem - (nElem & 7);
    for(int32_t j = 0; j < smLen; j+=8) {
      __m256i vIdx = _mm256_loadu_si256((__m256i*)&(indices[i][j]));
      __m256 vX = _mm256_i32gather_ps((float const*)x,vIdx,4);
      __m256 vA = _mm256_loadu_ps(&A[A_offset + j]);
      vT = _mm256_add_ps(vT, _mm256_mul_ps(vX,vA));
    }
    t += sum8(vT);
    for(int32_t j = smLen; j < nElem; j++) {
      int32_t idx = indices[i][j];
      t += x[idx]*A[A_offset + j];
    }
    y[i] = t;
    A_offset += nElem;
  }
}
```

Iterate over rows

For each 8 nonzeroes

Gather X, load A, multiply, add to y

Sum elements in vector

# SpMV for 8-wide SIMD

```
void avx2_csr_spmv( float *A, int32_t *nIdx, int32_t **indices, float *x, int32_t m, float *y) {
  int32_t A_offset = 0;
  for(int32_t i = 0; i < m; i++) {
    int32_t nElem = nIdx[i]; float t = 0.0f;
    __m256 vT = _mm256_set_ps(0.0f,0.0f,0.0f,0.0f,0.0f,0.0f,0.0f,0.0f);
    int32_t smLen = nElem - (nElem & 7);
    for(int32_t j = 0; j < smLen; j+=8) {
      __m256i vIdx = _mm256_loadu_si256((__m256i*)&(indices[i][j]));
      __m256 vX = _mm256_i32gather_ps((float const*)x,vIdx,4);
      __m256 vA = _mm256_loadu_ps(&A[A_offset + j]);
      vT = _mm256_add_ps(vT, _mm256_mul_ps(vX,vA));
    }
    t += sum8(vT);
    for(int32_t j = smLen; j < nElem; j++) {
      int32_t idx = indices[i][j];
      t += x[idx]*A[A_offset + j];
    }
    y[i] = t;
    A_offset += nElem;
  }
}
```

Iterate over rows

For each 8 nonzeroes

Gather X, load A, multiply, add to y

Sum elements in vector

In case #cols does not evenly divide vector width

# SpMV with CSR in CUDA

```
// Parallel SpMV using CSR format
__global__ void spmv(int *ptr, int *ind, float *val,
                int m, float *x, float *y) {
  for (int i = blockIdx.x * blockDim.x + threadIdx.x;
                i < m; i += blockDim.x * gridDim.x) {
    float yi = 0;
    for (int j = ptr[i]; j < ptr[i+1]; j++) {
      yi += values[j] * x[col_ind[j]];
    }
    y[i] = yi;
  }
}
```

Iterate over rows in A

# SpMV with CSR in CUDA

```
// Parallel SpMV using CSR format
__global__ void spmv(int *ptr, int *ind, float *val,
                int m, float *x, float *y) {
  for (int i = blockIdx.x * blockDim.x + threadIdx.x;
                i < m; i += blockDim.x * gridDim.x) {
    float yi = 0;
    for (int j = ptr[i]; j < ptr[i+1]; j++) {
      yi += values[j] * x[col_ind[j]];
    }
    y[i] = yi;
  }
}
```

Iterate over rows in A

Iterate over nonzeroes in row

# SpMV with CSR in CUDA

```
// Parallel SpMV using CSR format
__global__ void spmv(int *ptr, int *ind, float *val,
              int m, float *x, float *y) {
  for (int i = blockIdx.x * blockDim.x + threadIdx.x;
              i < m; i += blockDim.x * gridDim.x) {
    float yi = 0;
    for (int j = ptr[i]; j < ptr[i+1]; j++) {
      yi += values[j] * x[col_ind[j]];
    }
    y[i] = yi;
  }
}
```

**Iterate over rows in A**

**Iterate over nonzeroes in row**

**Assign results**

# Segmented Suffix Scan

|  | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|---|
| **x =** | 1 | 2 | 1 | 2 | 1 | 2 | 1 |

**ptr =** | 0 | 2 | 5 | 7 | 9 |

**A**

**ind =** | 1 | 2 | 0 | 3 | 4 | 1 | 2 | 5 | 6 | 0 | 1 | 4 |

**val =** | 1 | 1 | 1 | 2 | 2 | 2 | 1 | 2 | 2 | 1 | 3 | 1 |

flag = zeros + ones(ptr)

**flag =** | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 |

from x

**x(ind) =** | 2 | 1 | 1 | 2 | 1 | 2 | 1 | 2 | 1 | 1 | 2 | 1 |

prod = val * x(ind)

**prod =** | 2 | 1 | 1 | 4 | 2 | 4 | 1 | 4 | 2 | 1 | 6 | 1 |

# Segmented Suffix Scan

|  | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|---|
| x = | 1 | 2 | 1 | 2 | 1 | 2 | 1 |

| ptr = | 0 | 2 | 5 | 7 | 9 |
|---|---|---|---|---|---|

**A**

| ind = | 1 | 2 | 0 | 3 | 4 | 1 | 2 | 5 | 6 | 0 | 1 | 4 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

| val = | 1 | 1 | 1 | 2 | 2 | 2 | 1 | 2 | 2 | 1 | 3 | 1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

> **flag = zeros + ones(ptr)**

| flag = | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

> **prod = val * x(ind)**

| prod = | 2 | 1 | 1 | 4 | 2 | 4 | 1 | 4 | 2 | 1 | 6 | 1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

> **y = end of segments**

> **segmented scan on flag, prod**

| sums = | 2 | 3 | 1 | 5 | 7 | 4 | 5 | 4 | 6 | 1 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

# SpMV Diagonal Format in OpenMP

**vals:**

**# rows**

**# diagonals**

**columnoffset**

| |
|---|
| 1 |
| 0 |
| -1 |

```
for each diagonal k do
  #pragma omp parallel for
  for each row i do
    column = i + columnoffset[k]
    if (column >= 0 && column < n)
      y[i] = y[i] + val[k][column] * x[column]
```

# SpMV Diagonal Format in OpenMP

**vals:**

**# rows**

**# diagonals**

| 1 |
|---|
| 0 |
| -1 |

**columnoffset**

```
for each diagonal k do
  #pragma omp parallel for
  for each row i do
    column = i + columnoffset[k]
    if (column >= 0 && column < n)
      y[i] = y[i] + val[k][column] * x[column]
```

**Diagonal is also popular for GPUs/vectors**

# Distributed SpMV

# Distributed Dense Matrix-Vector Product (Row Major)

warmup on dense case

- Compute y = y + A*x, where A is a dense matrix
- Layout: **1D row blocked**
- Algorithm:

```
foreach processor p
  Broadcast x[p] chunk owned by p

for all local i
  for all j
    compute y[i] += A[i, j]*x[j] locally
```

**Broadcast + local dot product**

x

| | P0 |
|---|---|
| A(0) | |
| A(1) | P1 |
| A(2) | P2 |
| A(3) | P3 |

y

# Distributed Dense Matrix-Vector Product (Column Major)

**P0 P1 P2 P3**

**x**

- Compute y = y + A*x, where A is a dense matrix
- Layout: **1D column blocked**
- Algorithm:

**y**

A(0) A(1) A(2) A(3)

```
foreach processor p
  make a temp vector of size n
  for all local j
    for all i
      compute temp[i] += A[i, j]*x[j] locally

  y[i] += SumReduce(temp[i])
```

**Local daxpy and parallel reduction**

**Reduce across all rows**

# Distributed Dense Matrix-Vector Product (2D Blocked)

A 2D blocked layout uses a broadcast of x and reduction into y.

Both on a subset of processors
- sqrt(p) for square processor grid
- Can use other rectangular shapes p_row * p_col = p



**x**

**y**

| P(0) | P(1) | P(2) | P(3) |
| P(4) | P(5) | P(6) | P(7) |
| P(8) | P(9) | P(10) | P(11) |
| P(12) | P(13) | P(14) | P(15) |

# Distributed SpMV

Row parallelism (y & A partitioned)
- Replicate x across processors
- Or exchange only necessary elements
- Are nonzeroes clustered, e.g., near the diagonal?

# Distributed SpMV

Column parallelism (x & A partitioned)
- Make temporary temp_y = [0,…] on all processors;
- Update that; and (sparse?) sum-reduce over processors

**P0  P1  P2  P3**

**x**

**y**

A(0) A(1) A(2) A(3)

# Distributed SpMV

2D parallelism for large p and **when nonzeros are uniform**
- Divide processors into p1 x p2 (e.g., square grid)
- Hybrid of Row and Column parallelism using teams
- NAS CG benchmark does this (random nonzero pattern)
- Bad load balance for clustered nonzeros

# Ideal Sparse Structure: P Diagonal Blocks

"Ideal" matrix structure for parallelism: block diagonal

- If $p_i$ holds $x_i$ and $y_i$ blocks, no vectors to communicate
- If non nonzeroes outside these blocks, no communication is needed!



Dream scenario: reorder rows/columns to get close to this - most nonzeroes in diagonal blocks, few outside.

# High Performance Conjugate Gradients (HPCG) Benchmark

**Complement to LINPACK** (dense linear algebra), which is used for the TOP500.

Designed to exercise computational and data access patterns that more closely match a **different and broader set of applications**.

HPCG includes performance of the following basic operations:
- Conjugate gradient (27-point stencil)
- Sparse matrix-vector multiplication
- Global dot product
- Vector update
- and others



$$\mathcal{L}[u] \equiv \nabla^2 u = f$$

$$f(\Omega) = 0$$

27-point discretization of a regular 3-D grid

GEOEMTRIC MULTIGRID PRECONDITIONER

3-D Data Grid mapped onto a 3-D Process Grid of MPI

Halo Exchange via MPI

$$Au = f$$

## November 2023 HPCG Results

**New HPCG results were announced at SC23**

| Rank | Site | Computer | Cores | HPL Rmax (Pflop/s) | TOP500 Rank | HPCG (Pflop/s) | Fraction of Peak |
|---|---|---|---|---|---|---|---|
| 1 | RIKEN Center for Computational Science **Japan** | **Supercomputer Fugaku** — A64FX 48C 2.2GHz, Tofu interconnect D | 7,630,848 | 442.01 | 4 | 16.00 | 3.0% |
| 2 | DOE/SC/Oak Ridge National Laboratory **United States** | **Frontier** — AMD Optimized 3rd Generation EPYC 64C 2GHz, Slingshot-11, AMD Instinct MI250X | 8,699,904 | 1194.00 | 1 | 14.05 | 0.8% |
| 3 | EuroHPC/CSC **Finland** | **LUMI** — AMD Optimized 3rd Generation EPYC 64C 2GHz, Slingshot-11, AMD Instinct MI250X | 2,752,704 | 379.70 | 5 | 4.587 | 0.9% |
| 4 | EuroHPC/CINECA **Italy** | **Leonardo** — Xeon Platinum 8358 32C 2.6GHz, Quad-rail NVIDIA HDR100 Infiniband, NVIDIA A100 SXM4 64 GB | 1,824,768 | 238.70 | 6 | 3.114 | 1.0% |
| 5 | DOE/SC/Oak Ridge National Laboratory **United States** | **Summit** — IBM POWER9 22C 3.07GHz, Dual-rail Mellanox EDR Infiniband, NVIDIA Volta GV100 | 2,414,592 | 148.60 | 7 | 2.926 | 1.5% |
| 6 | DOE/SC/LBNL/NERSC **United States** | **Perlmutter** — AMD EPYC 7763 64C 2.45GHz, Slingshot-11, NVIDIA A100 SXM4 40 GB | 888,832 | 79.23 | 12 | 1.905 | 1.7% |
| 7 | DOE/NNSA/LLNL **United States** | **Sierra** — IBM POWER9 22C 3.1GHz, Dual-rail Mellanox EDR Infiniband, NVIDIA Volta GV100 | 1,572,480 | 94.64 | 10 | 1.796 | 1.4% |
| 8 | NVIDIA Corporation **United States** | **Selene** — AMD EPYC 7742 64C 2.25GHz, Mellanox HDR Infiniband, NVIDIA A100 | 555,520 | 63.46 | 13 | 1.623 | 2.0% |
| 9 | Forschungszentrum Juelich (FZJ) **Germany** | **JUWELS Booster Module** — AMD EPYC 7402 24C 2.8GHz, Mellanox HDR InfiniBand/ParTec ParaStation ClusterSuite, NVIDIA A100 | 449,280 | 44.12 | 18 | 1.275 | 1.8% |
| 10 | Cyberscience Center, Tohoku University **Japan** | **AOBA-S** — Vector Engine Type 30A 16C 1.6GHz, Infiniband NDR200 | 64,512 | 17.22 | 50 | 1.089 | 5.5% |

Not the same order as TOP500

## November 2023 HPCG Results

https://www.netlib.org/benchmark/hpl/

| Rank | Site | Computer | Cores | HPL Rmax (Pflop/s) | TOP500 Rank | HPCG (Pflop/s) | Fraction of Peak |
|------|------|----------|-------|--------------------|-------------|-----------------|------------------|
| 1 | RIKEN Center for Computational Science **Japan** | **Supercomputer Fugaku** — A64FX 48C 2.2GHz, Tofu interconnect D | 7,630,848 | 442.01 | 4 | 16.00 | 3.0% |
| 2 | DOE/SC/Oak Ridge National Laboratory **United States** | **Frontier** — AMD Optimized 3rd Generation EPYC 64C 2GHz, Slingshot-11, AMD Instinct MI250X | 8,699,904 | 1194.00 | 1 | 14.05 | 0.8% |
| 3 | EuroHPC/CSC **Finland** | **LUMI** — AMD Optimized 3rd Generation EPYC 64C 2GHz, Slingshot-11, AMD Instinct MI250X | 2,752,704 | 379.70 | 5 | 4.587 | 0.9% |
| 4 | EuroHPC/CINECA **Italy** | **Leonardo** — Xeon Platinum 8358 32C 2.6GHz, Quad-rail NVIDIA HDR100 Infiniband, NVIDIA A100 SXM4 64 GB | 1,824,768 | 238.70 | 6 | 3.114 | 1.0% |
| 5 | DOE/SC/Oak Ridge National Laboratory **United States** | **Summit** — IBM POWER9 22C 3.07GHz, Dual-rail Mellanox EDR Infiniband, NVIDIA Volta GV100 | 2,414,592 | 148.60 | 7 | 2.926 | 1.5% |
| 6 | DOE/SC/LBNL/NERSC **United States** | **Perlmutter** — AMD EPYC 7763 64C 2.45GHz, Slingshot-11, NVIDIA A100 SXM4 40 GB | 888,832 | 79.23 | 12 | 1.905 | 1.7% |
| 7 | DOE/NNSA/LLNL **United States** | **Sierra** — IBM POWER9 22C 3.1GHz, Dual-rail Mellanox EDR Infiniband, NVIDIA Volta GV100 | 1,572,480 | 94.64 | 10 | 1.796 | 1.4% |
| 8 | NVIDIA Corporation **United States** | **Selene** — AMD EPYC 7742 64C 2.25GHz, Mellanox HDR Infiniband, NVIDIA A100 | 555,520 | 63.46 | 13 | 1.623 | 2.0% |
| 9 | Forschungszentrum Juelich (FZJ) **Germany** | **JUWELS Booster Module** — AMD EPYC 7402 24C 2.8GHz, Mellanox HDR InfiniBand/ParTec ParaStation ClusterSuite, NVIDIA A100 | 449,280 | 44.12 | 18 | 1.275 | 1.8% |
| 10 | Cyberscience Center, Tohoku University **Japan** | **AOBA-S** — Vector Engine Type 30A 16C 1.6GHz, Infiniband NDR200 | 64,512 | 17.22 | 50 | 1.089 | 5.5% |

**Much lower Pflop/s compared to Linpack**

## November 2023 HPCG Results
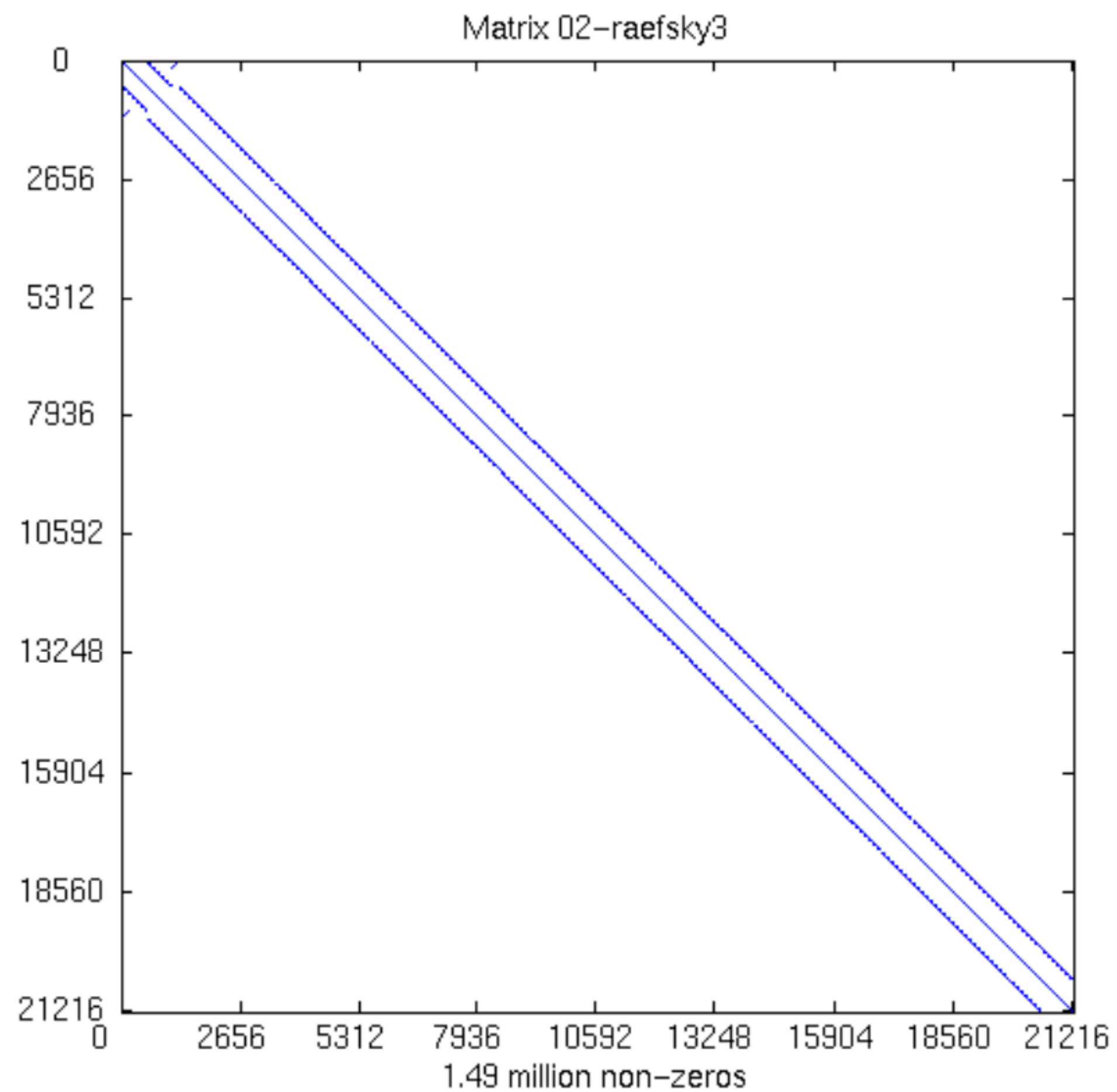
**New HPCG results were announced at SC23**

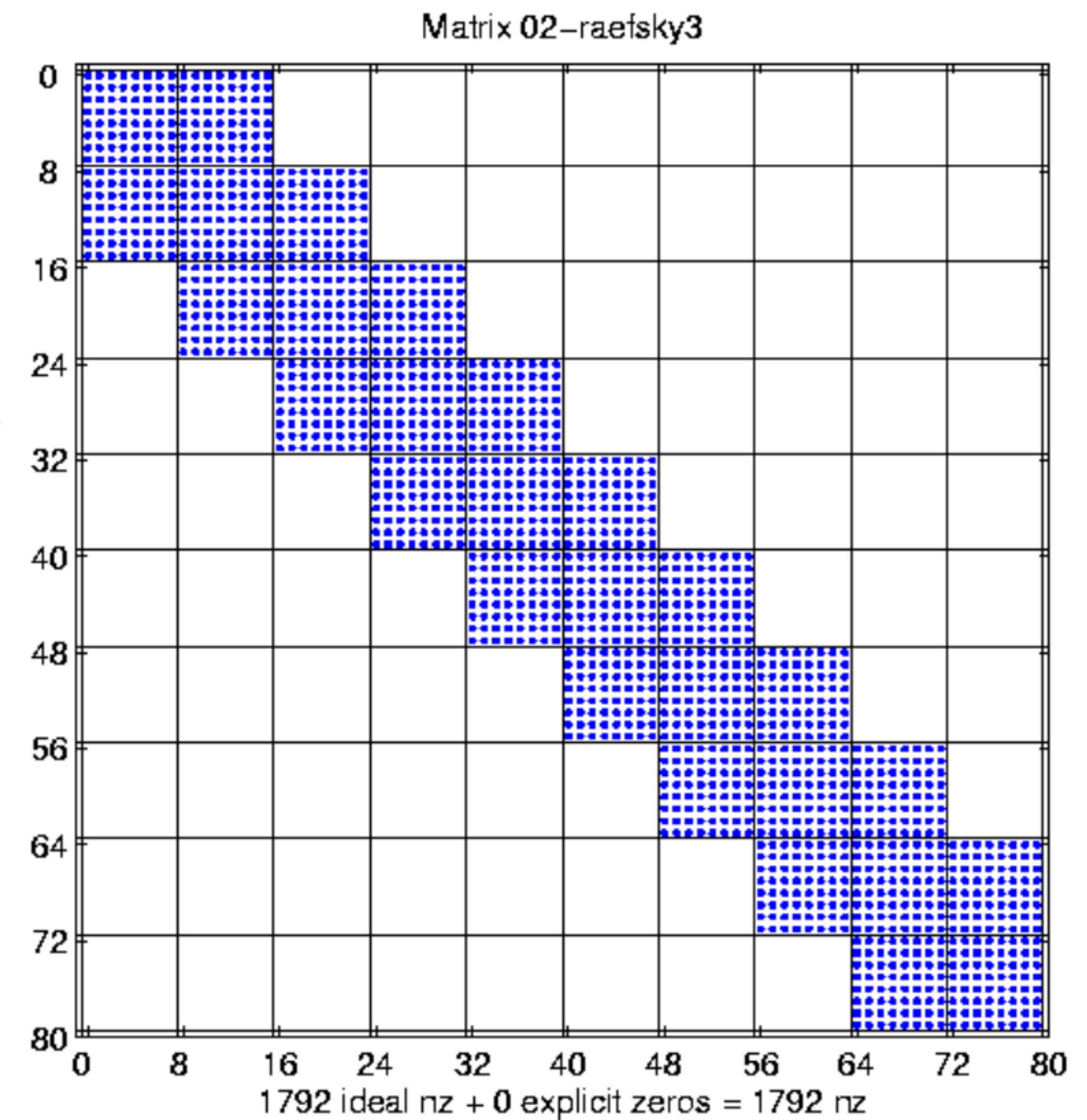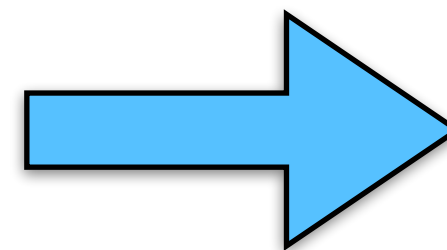| Rank | Site | Computer | Cores | HPL Rmax (Pflop/s) | TOP500 Rank | HPCG (Pflop/s) | Fraction of Peak |
|------|------|----------|-------|--------------------|-------------|----------------|------------------|
| 1 | RIKEN Center for Computational Science **Japan** | **Supercomputer Fugaku** — A64FX 48C 2.2GHz, Tofu interconnect D | 7,630,848 | 442.01 | 4 | 16.00 | 3.0% |
| 2 | DOE/SC/Oak Ridge National Laboratory **United States** | **Frontier** — AMD Optimized 3rd Generation EPYC 64C 2GHz, Slingshot-11, AMD Instinct MI250X | 8,699,904 | 1194.00 | 1 | 14.05 | 0.8% |
| 3 | EuroHPC/CSC **Finland** | **LUMI** — AMD Optimized 3rd Generation EPYC 64C 2GHz, Slingshot-11, AMD Instinct MI250X | 2,752,704 | 379.70 | 5 | 4.587 | 0.9% |
| 4 | EuroHPC/CINECA **Italy** | **Leonardo** — Xeon Platinum 8358 32C 2.6GHz, Quad-rail NVIDIA HDR100 Infiniband, NVIDIA A100 SXM4 64 GB | 1,824,768 | 238.70 | 6 | 3.114 | 1.0% |
| 5 | DOE/SC/Oak Ridge National Laboratory **United States** | **Summit** — IBM POWER9 22C 3.07GHz, Dual-rail Mellanox EDR Infiniband, NVIDIA Volta GV100 | 2,414,592 | 148.60 | 7 | 2.926 | 1.5% |
| 6 | DOE/SC/LBNL/NERSC **United States** | **Perlmutter** — AMD EPYC 7763 64C 2.45GHz, Slingshot-11, NVIDIA A100 SXM4 40 GB | 888,832 | 79.23 | 12 | 1.905 | 1.7% |
| 7 | DOE/NNSA/LLNL **United States** | **Sierra** — IBM POWER9 22C 3.1GHz, Dual-rail Mellanox EDR Infiniband, NVIDIA Volta GV100 | 1,572,480 | 94.64 | 10 | 1.796 | 1.4% |
| 8 | NVIDIA Corporation **United States** | **Selene** — AMD EPYC 7742 64C 2.25GHz, Mellanox HDR Infiniband, NVIDIA A100 | 555,520 | 63.46 | 13 | 1.623 | 2.0% |
| 9 | Forschungszentrum Juelich (FZJ) **Germany** | **JUWELS Booster Module** — AMD EPYC 7402 24C 2.8GHz, Mellanox HDR InfiniBand/ParTec ParaStation ClusterSuite, NVIDIA A100 | 449,280 | 44.12 | 18 | 1.275 | 1.8% |
| 10 | Cyberscience Center, Tohoku University **Japan** | **AOBA-S** — Vector Engine Type 30A 16C 1.6GHz, Infiniband NDR200 | 64,512 | 17.22 | 50 | 1.089 | 5.5% |

**Small fraction of peak**

# Register/cache blocking and autotuning SpMV

# Changing Matrix Format: Register Blocking



**Submatrix**

"Fast sparse matrix-vector multiplication by exploiting variable block structure," Vuduc and Moon 2005.
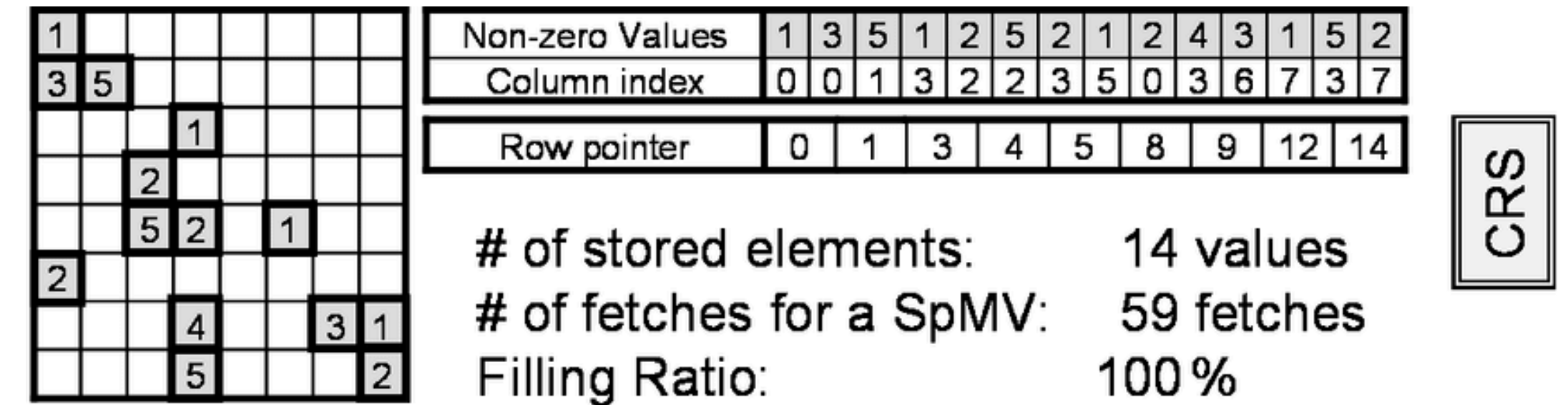
# Using block structure in SpMV

The **bottleneck is the time to fetch** the matrix from memory
- Only 2 flops for each nz in matrix
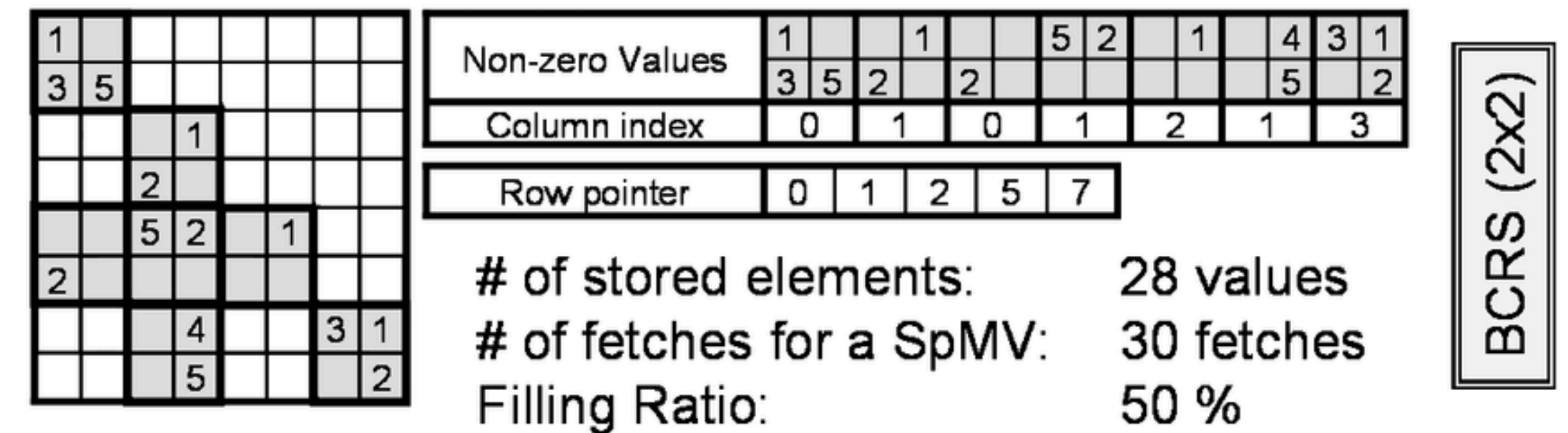- Fetching at ~1 int (col idx) + 1 float (value) for 2 flops

Blocked Compressed Sparse Row (BCSR)
- Don't store each nonzero - instead, **store each nonzero r x c block** with 1 column index
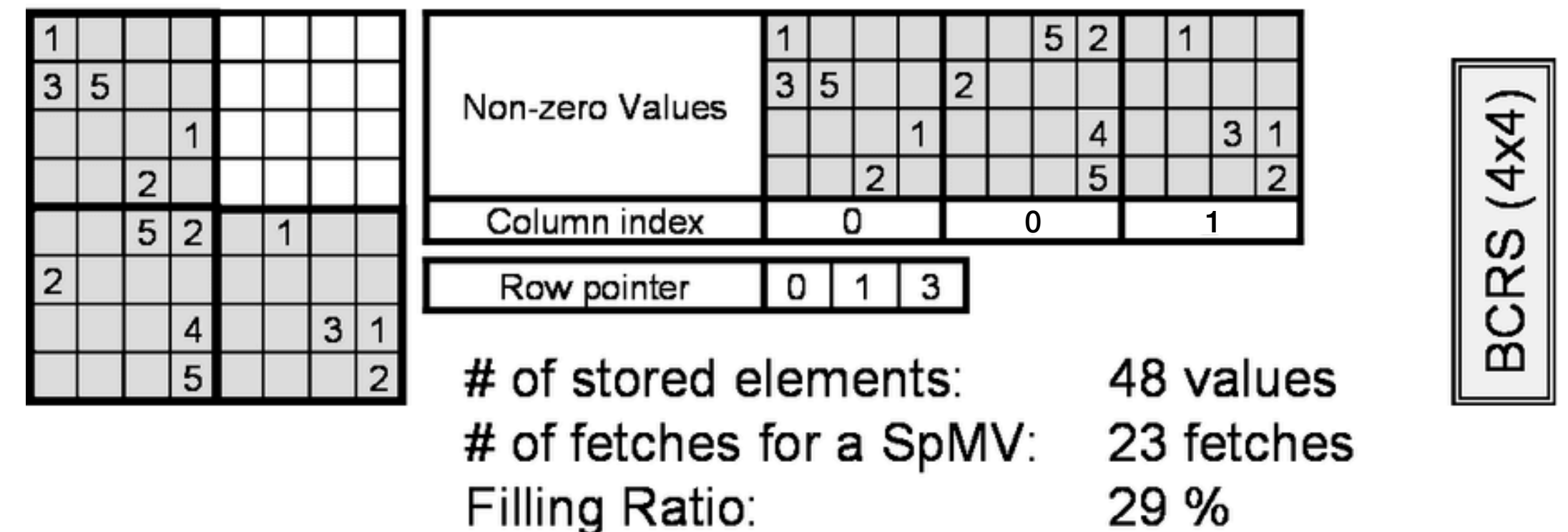- Time to fetch matrix from memory decreases

**Change both data structure and algorithm** - need to pick r and c and change algorithm accordingly



CRS

| Non-zero Values | 1 | 3 | 5 | 1 | 2 | 5 | 2 | 1 | 2 | 4 | 3 | 1 | 5 | 2 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Column index | 0 | 0 | 1 | 3 | 2 | 2 | 3 | 5 | 0 | 3 | 6 | 7 | 3 | 7 |

| Row pointer | 0 | 1 | 3 | 4 | 5 | 8 | 9 | 12 | 14 |
|---|---|---|---|---|---|---|---|---|---|

# of stored elements:        14 values
# of fetches for a SpMV:     59 fetches
Filling Ratio:               100 %

BCRS (2x2)

| Non-zero Values | 1 | | | 1 | | 5 | 2 | | 1 | | 4 | 3 | 1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 3 | 5 | 2 | | 2 | | | | | | 5 | | 2 |

| Column index | 0 | 1 | 0 | 1 | 2 | 1 | 3 |
|---|---|---|---|---|---|---|---|

| Row pointer | 0 | 1 | 2 | 5 | 7 |
|---|---|---|---|---|---|

# of stored elements:        28 values
# of fetches for a SpMV:     30 fetches
Filling Ratio:               50 %

BCRS (4x4)

| Non-zero Values | 1 | | | | | 5 | 2 | | 1 | |
|---|---|---|---|---|---|---|---|---|---|---|
| | 3 | 5 | | 2 | | | | | | |
| | | | 1 | | | 4 | | 3 | 1 |
| | | 2 | | | 5 | | | 2 |

| Column index | 0 | 0 | 1 |
|---|---|---|---|

| Row pointer | 0 | 1 | 3 |
|---|---|---|---|

# of stored elements:        48 values
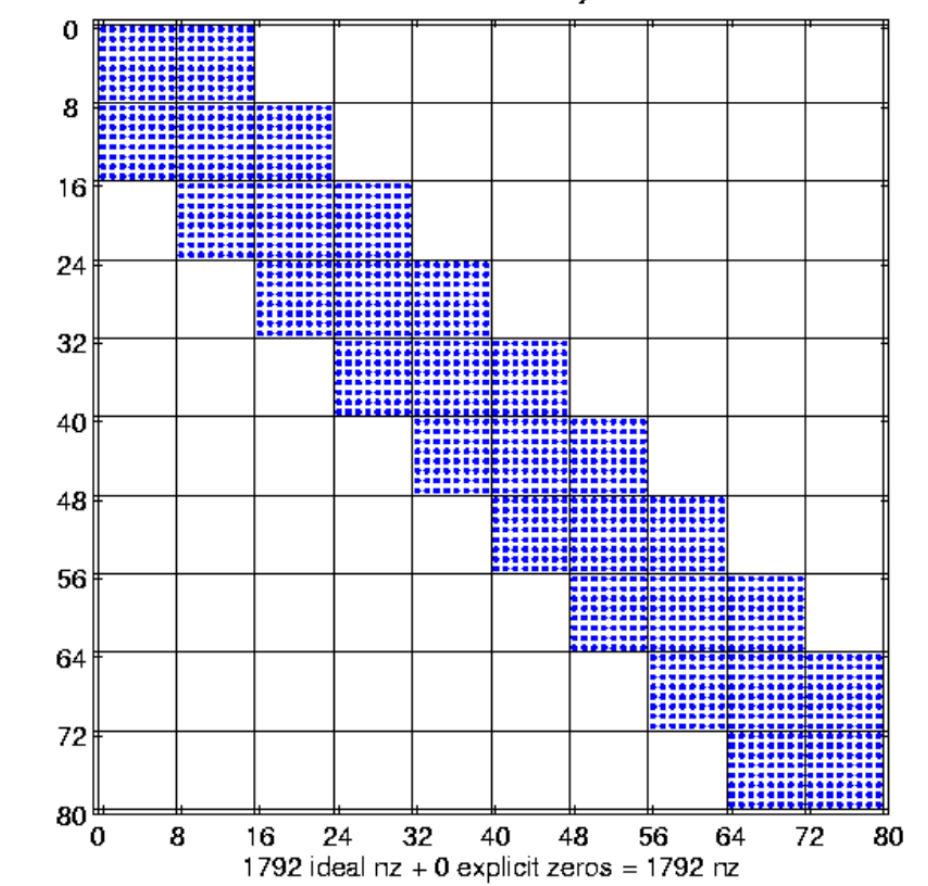# of fetches for a SpMV:     23 fetches
Filling Ratio:               29 %

# The need for search

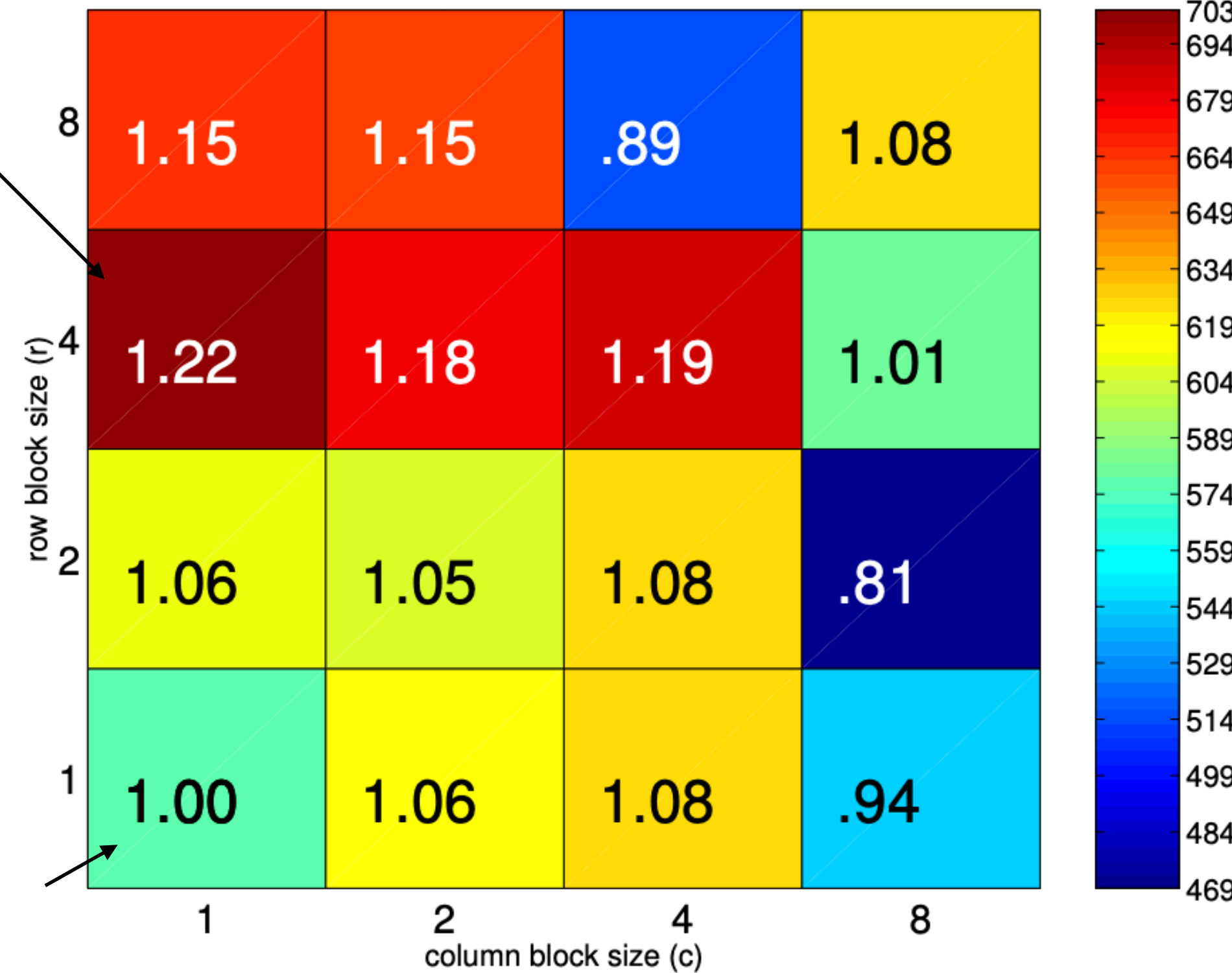In the previous example, it seems like 8x8 is a natural choice.

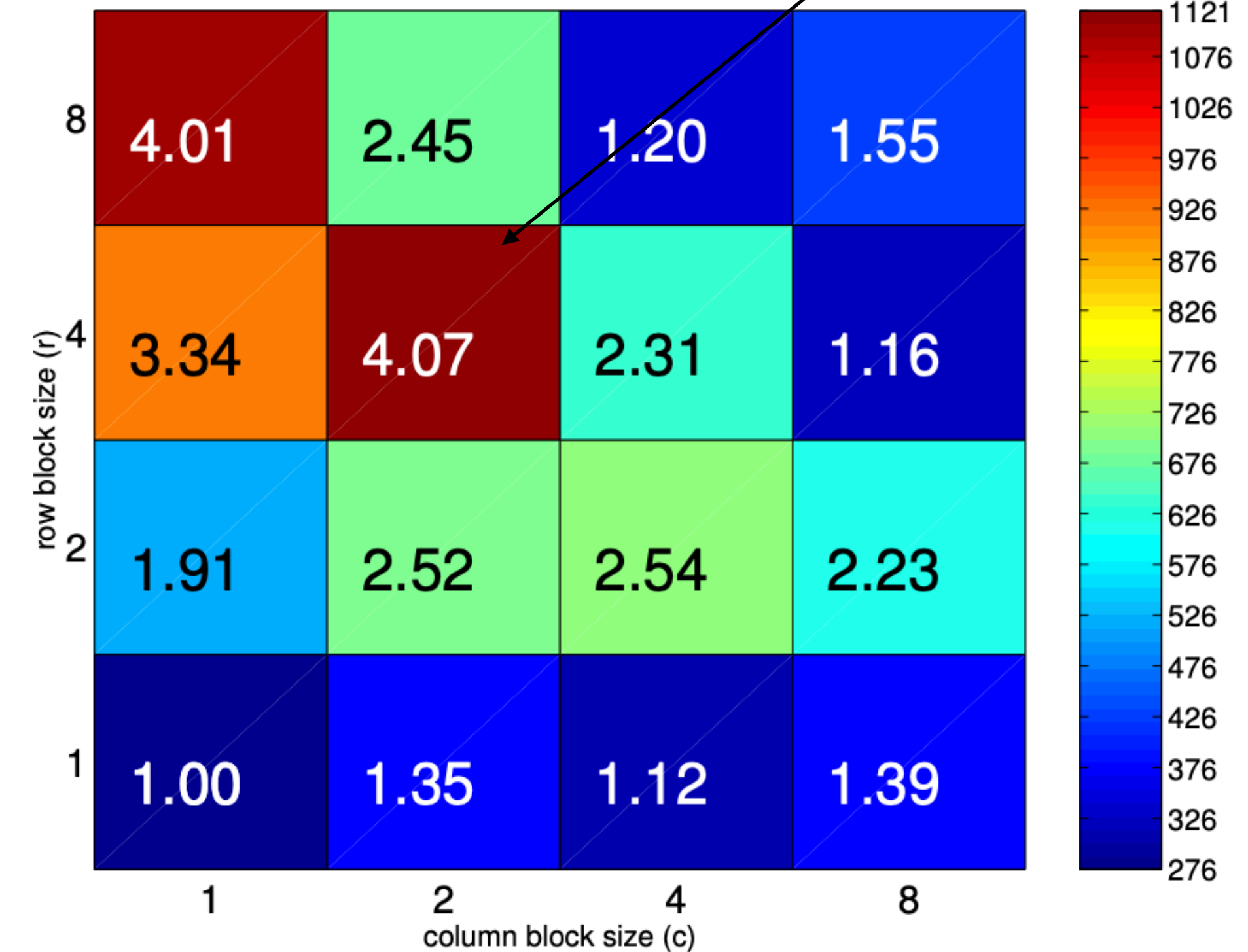Out of 6 platforms tested, 8x8 is the best on one (coming up)

Best: 4x2

Best: 4x1

Reference



"Automatic Performance Tuning of Sparse Matrix Kernels," Vuduc, 2003.
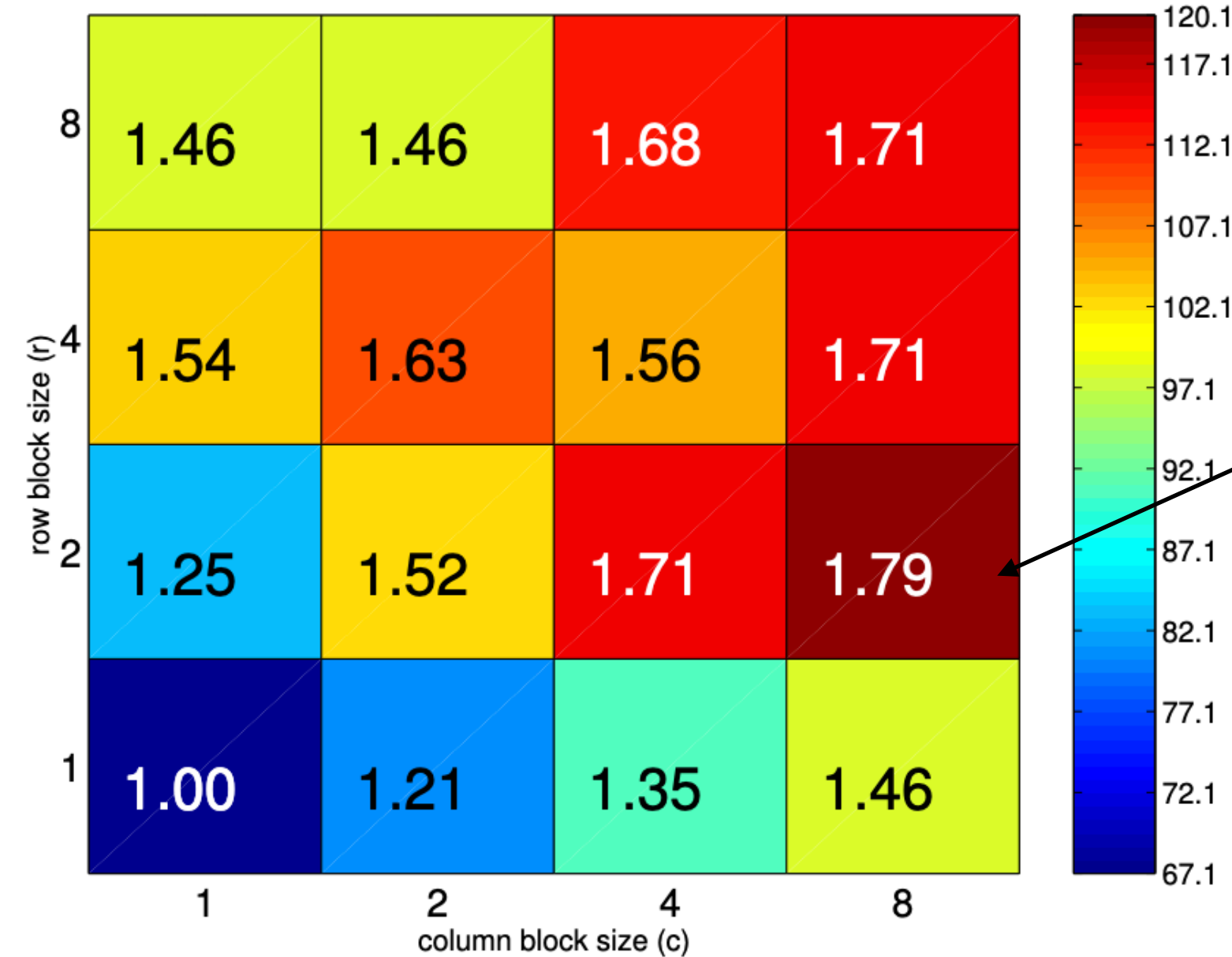
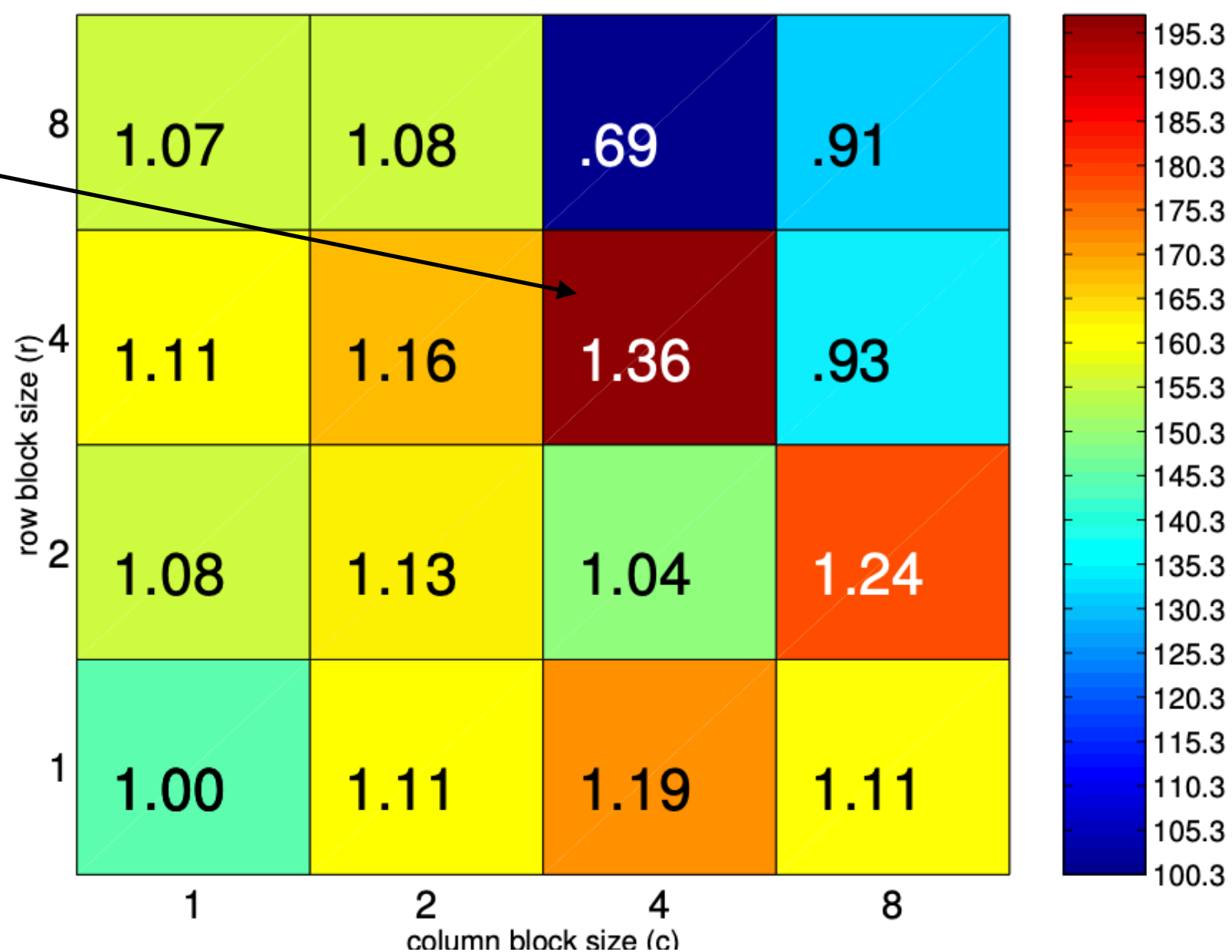SpMV Performance: raefsky3.rua [ref=35.3 Mflop/s; 333 MHz Sun Ultra 2i, Sun C v6.0]

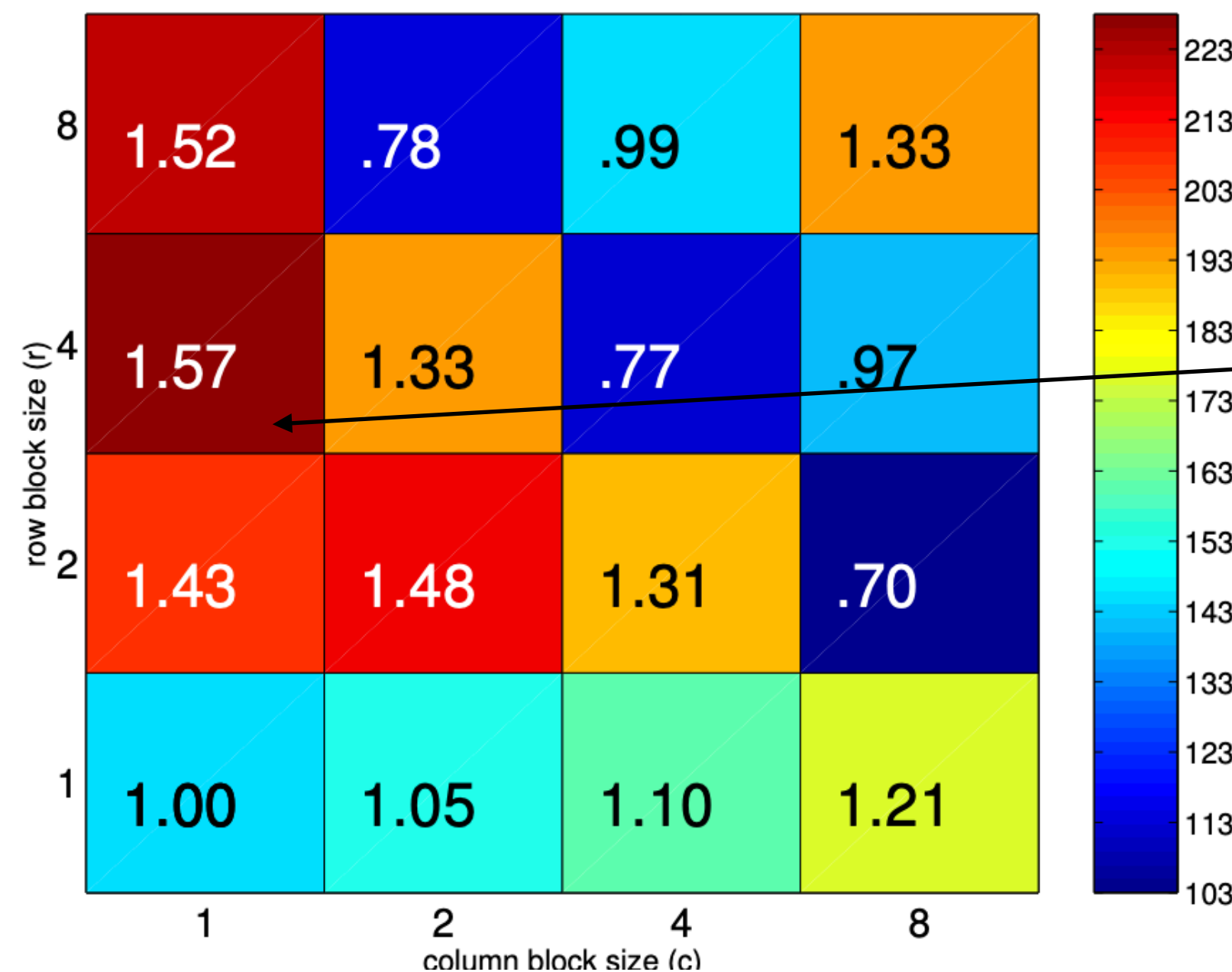SpMV Performance: raefsky3.rua [ref=67.1 Mflop/s; 800 MHz Pentium III–M, Intel C v7.0]

SpMV Performance: raefsky3.rua [ref=144.7 Mflop/s; 375 MHz Power3, IBM xlc v6]

SpMV Performance: raefsky3.rua [ref=145.8 Mflop/s; 800 MHz Itanium, Intel C v7]
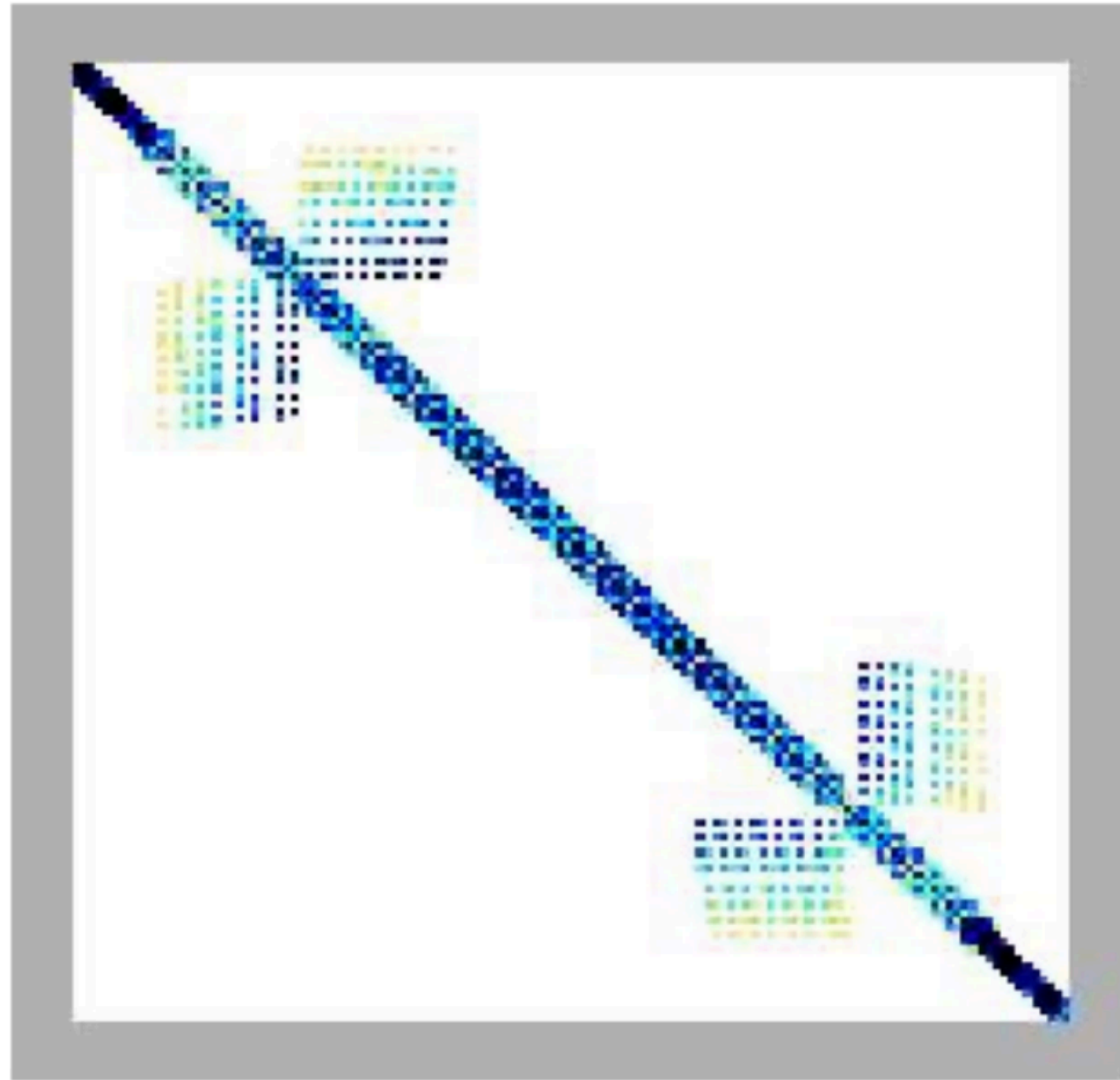
Best: 8x8

Best: 2x8
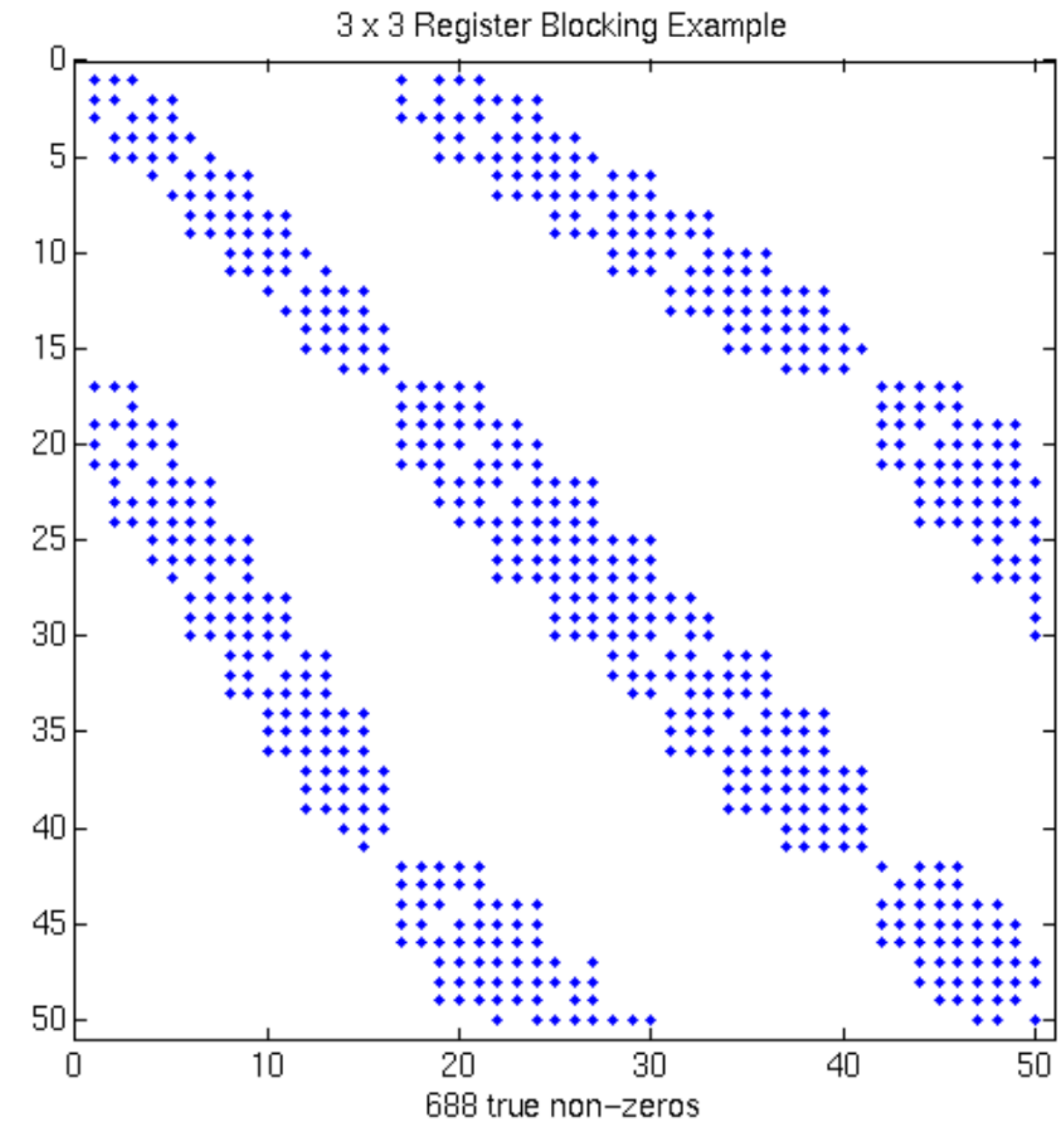
Best: 4x4

Best: 4x1

"Automatic Performance Tuning of Sparse Matrix Kernels," Vuduc, 2003.

# But most matrices don't block so easily
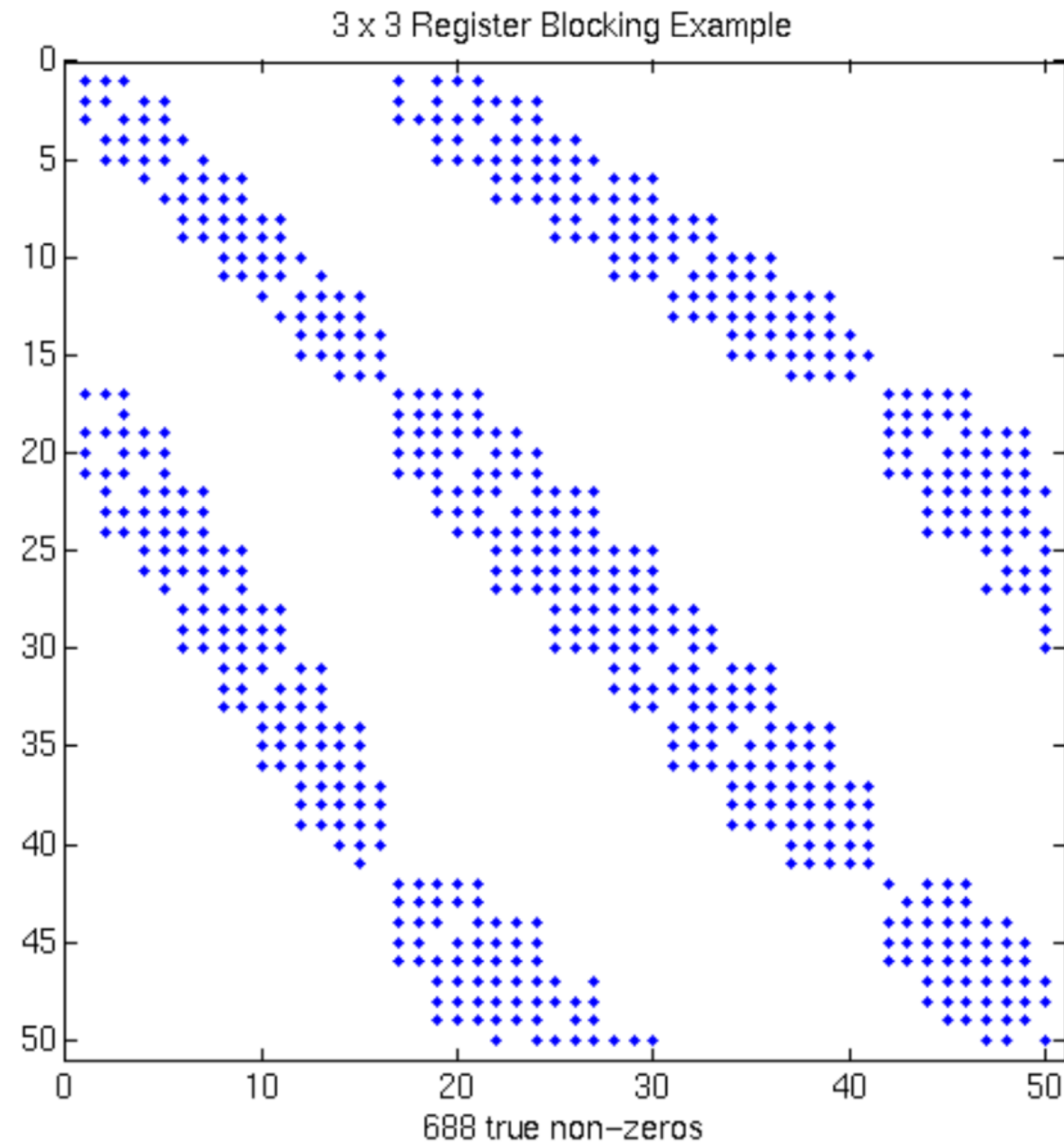


**Zoom into top corner** →

3 x 3 Register Blocking Example

688 true non-zeros

Fluid dynamics problem - more complicated nonzero structure
Total nnz = 1.1M

# 3x3 blocks look natural, but…

3x3 blocks

Many blocks are not full
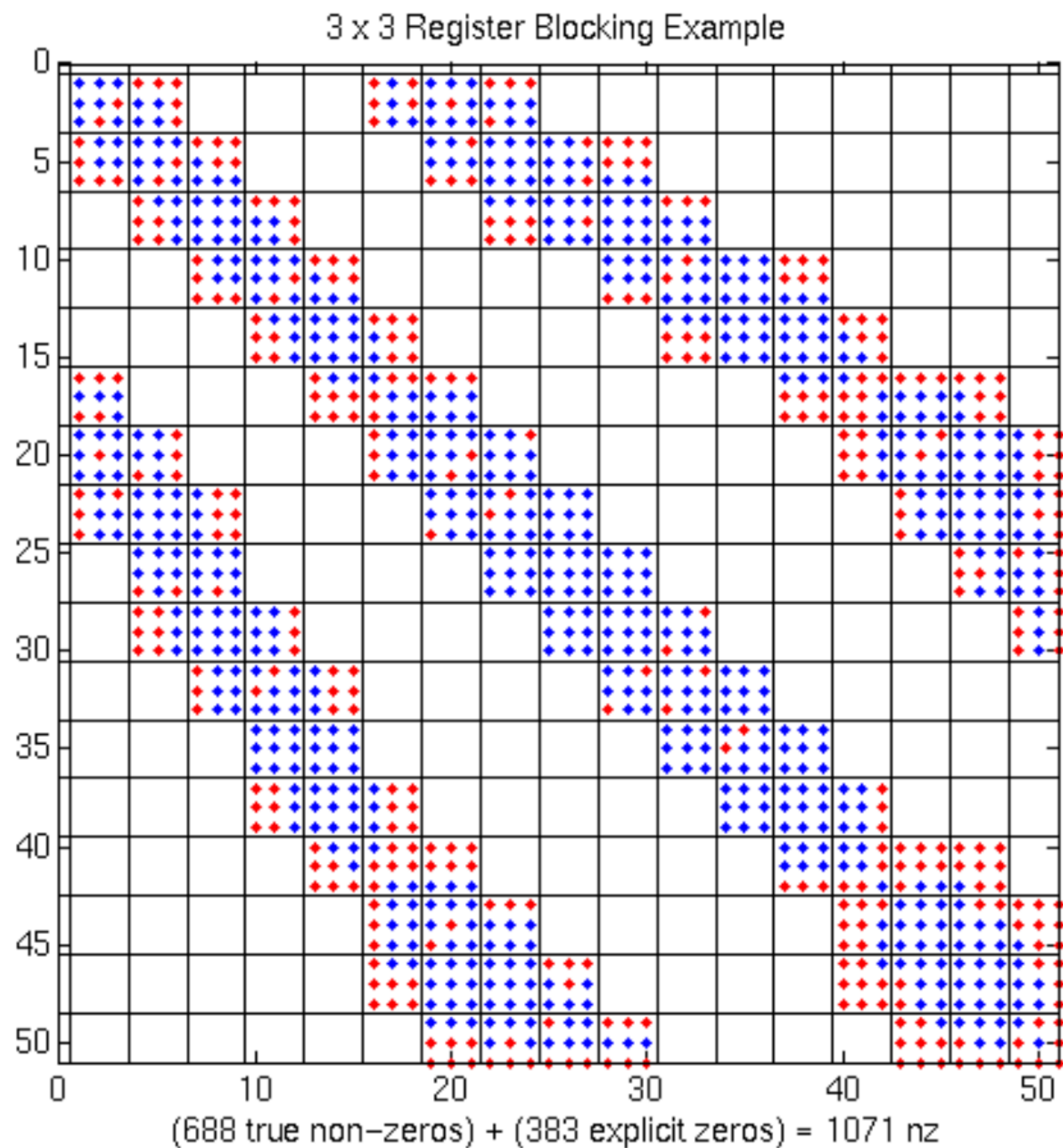
# Extra work can improve efficiency



3 x 3 Register Blocking Example

(688 true non-zeros) + (383 explicit zeros) = 1071 nz

- Add explicit zeroes: 1.5x "fill overhead"

- Unroll loops

- More work but faster - 1.5x faster on PIII

# Libraries for sparse matrices

How to build optimized library when:
- Formats are not known? Libraries like PETSc and Trilinos will let the user provide format and SpMV

How to build optimized matrix kernel library (BLAS)?
- Nonzero structure is key to optimization

OSKI = Optimized Sparse Kernel Interface
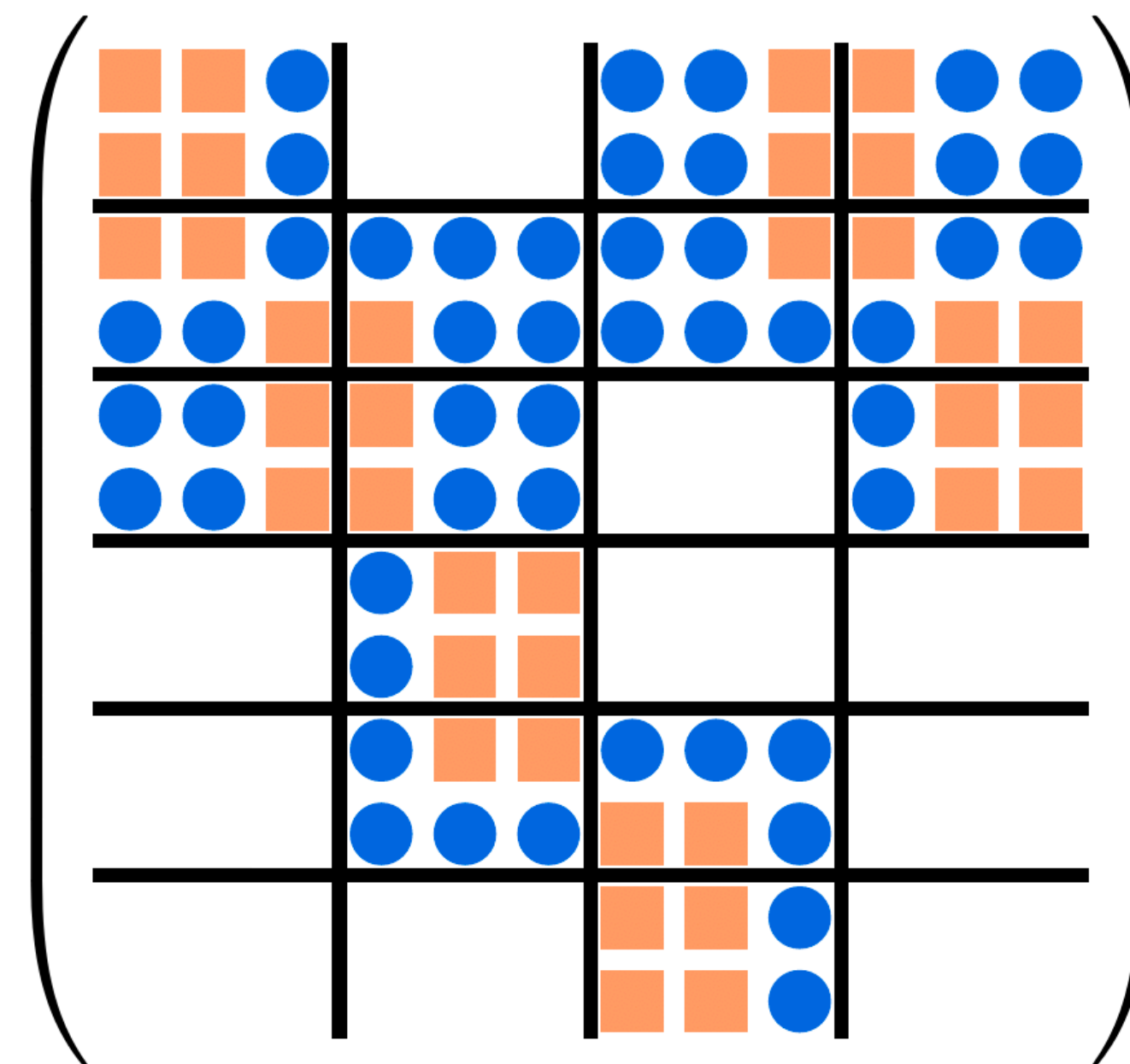- pOSKI for multicore

# Using the Fill to Formalize Blocking Scheme Quality [Im & Yelick, 01]

The fill of a *N*-dimensional tensor *A* is defined with respect to
- the number of nonzeros $k(A)$
- a blocking $\mathbf{b} = (b_1, \ldots, b_N)$
- the number of nonempty blocks $k_{\mathbf{b}}(A)$

$$f_{\mathbf{b}}(A) = \frac{b_1 b_2 \ldots b_N k_{\mathbf{b}}(A)}{k(A)}$$

# Using the Fill to Formalize Blocking Scheme Quality [Im & Yelick, 01]

The fill of a *N*-dimensional tensor *A* is defined with respect to
- the number of nonzeros $k(A)$
- a blocking $\mathbf{b} = (b_1, \ldots, b_N)$
- the number of nonempty blocks $k_{\mathbf{b}}(A)$
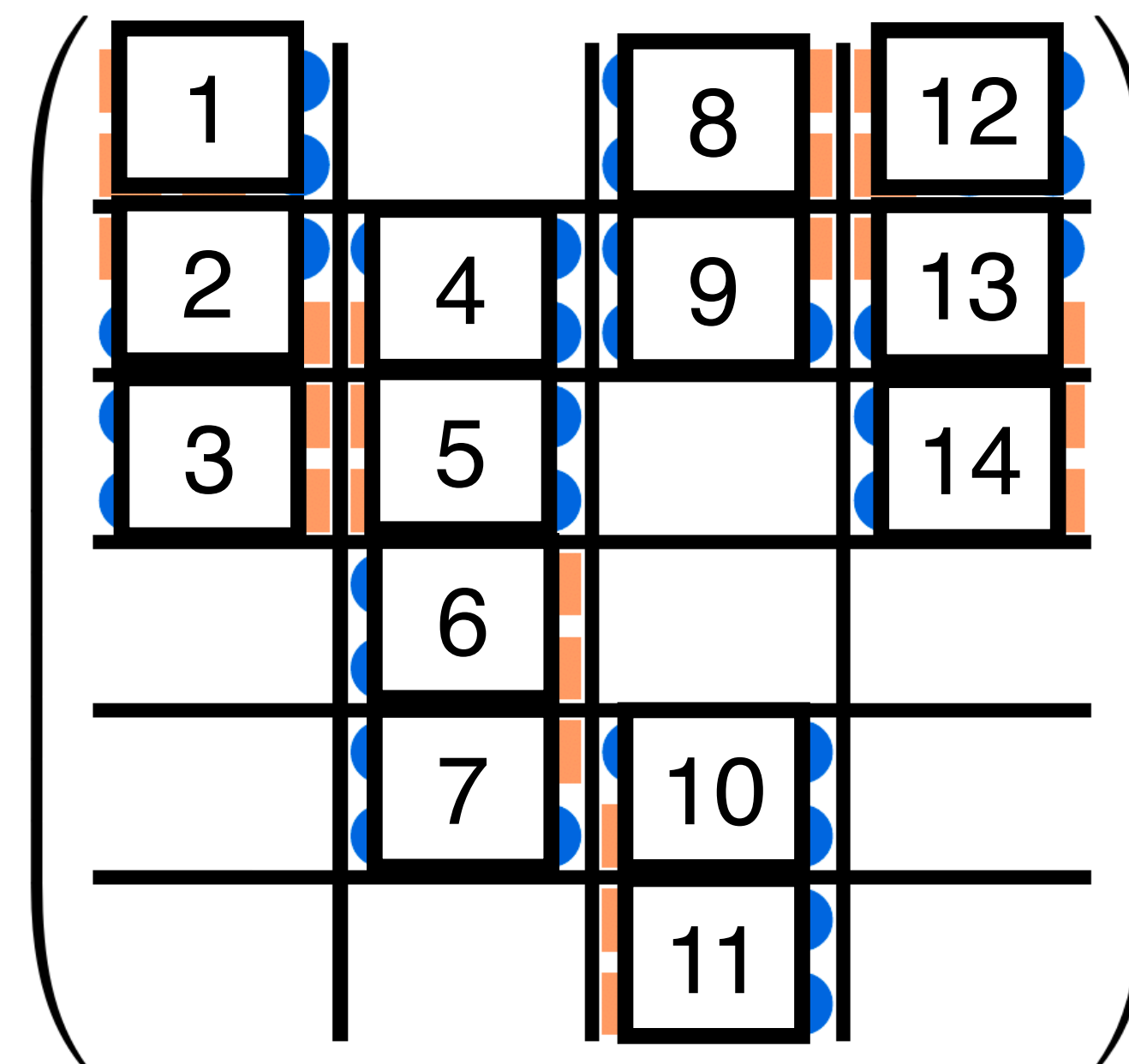
$$f_{\mathbf{b}}(A) = \frac{b_1 b_2 \ldots b_N k_{\mathbf{b}}(A)}{k(A)}$$

$$= \frac{2 \times 3 \times 14}{36} \approx 2.33$$

Nonempty blocks

Nonzeros

Greater ~ more wasted space
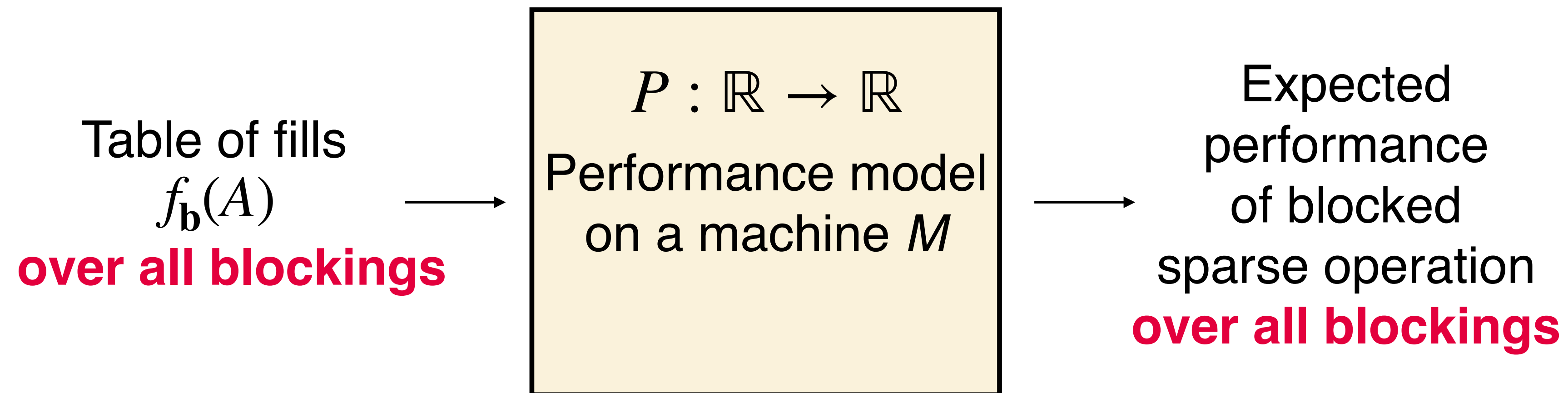
**Blocking** $\mathbf{b}$

$b_1 = 2$
$b_2 = 3$
$k_{\mathbf{b}} = 14$
$k(\mathcal{A}) = 36$



59

# Performance Modeling Using the Fill

A **performance model** is a function that maps the fill under a blocking $\mathbf{b}$ to expected performance in FLOP/s.

Table of fills
$f_{\mathbf{b}}(A)$
**over all blockings**

$\longrightarrow$

$$P : \mathbb{R} \to \mathbb{R}$$
Performance model on a machine *M*

$\longrightarrow$

Expected performance of blocked sparse operation
**over all blockings**

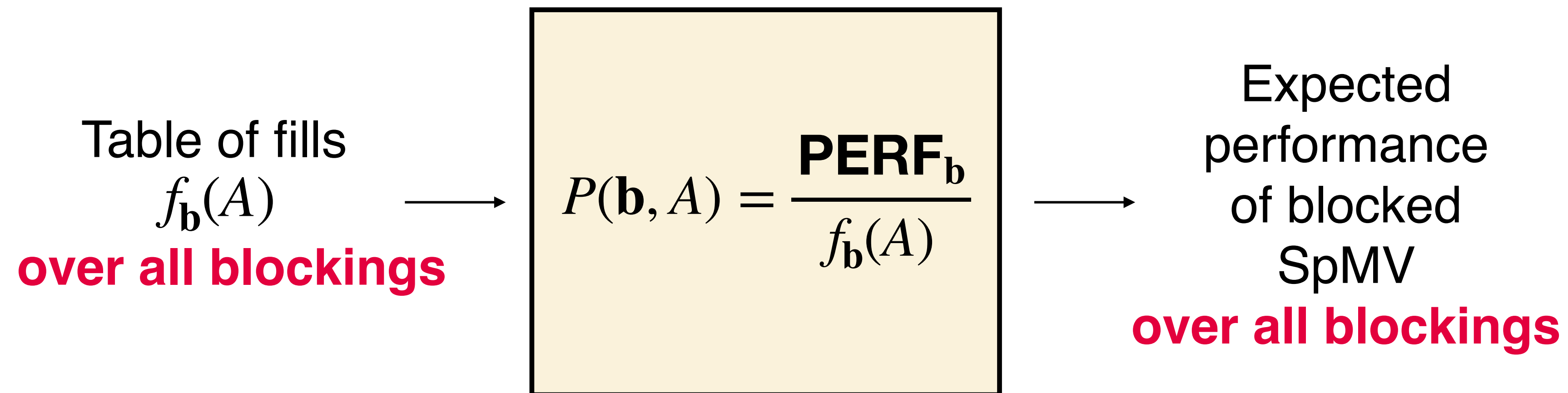**Block size selection** chooses the block size that maximizes performance based on a performance model.

"Optimizing the Performance of Sparse Matrix-Vector Multiplication," Im, 2000.
"Automatic Performance Tuning of Sparse Matrix Kernels," Vuduc, 2003.

# Example: SPARSITY Performance Model for Blocked SpMV [Vuduc et al., 02]

The SPARSITY performance model uses a matrix $\mathbf{PERF}(\mathbf{b})$ of the performance of a machine $M$ (in FLOP/s) on a dense matrix stored with blocking scheme $\mathbf{b}$.

Table of fills
$f_{\mathbf{b}}(A)$
**over all blockings**

$\longrightarrow$

$$P(\mathbf{b}, A) = \frac{\mathbf{PERF_b}}{f_{\mathbf{b}}(A)}$$

$\longrightarrow$

Expected performance of blocked SpMV
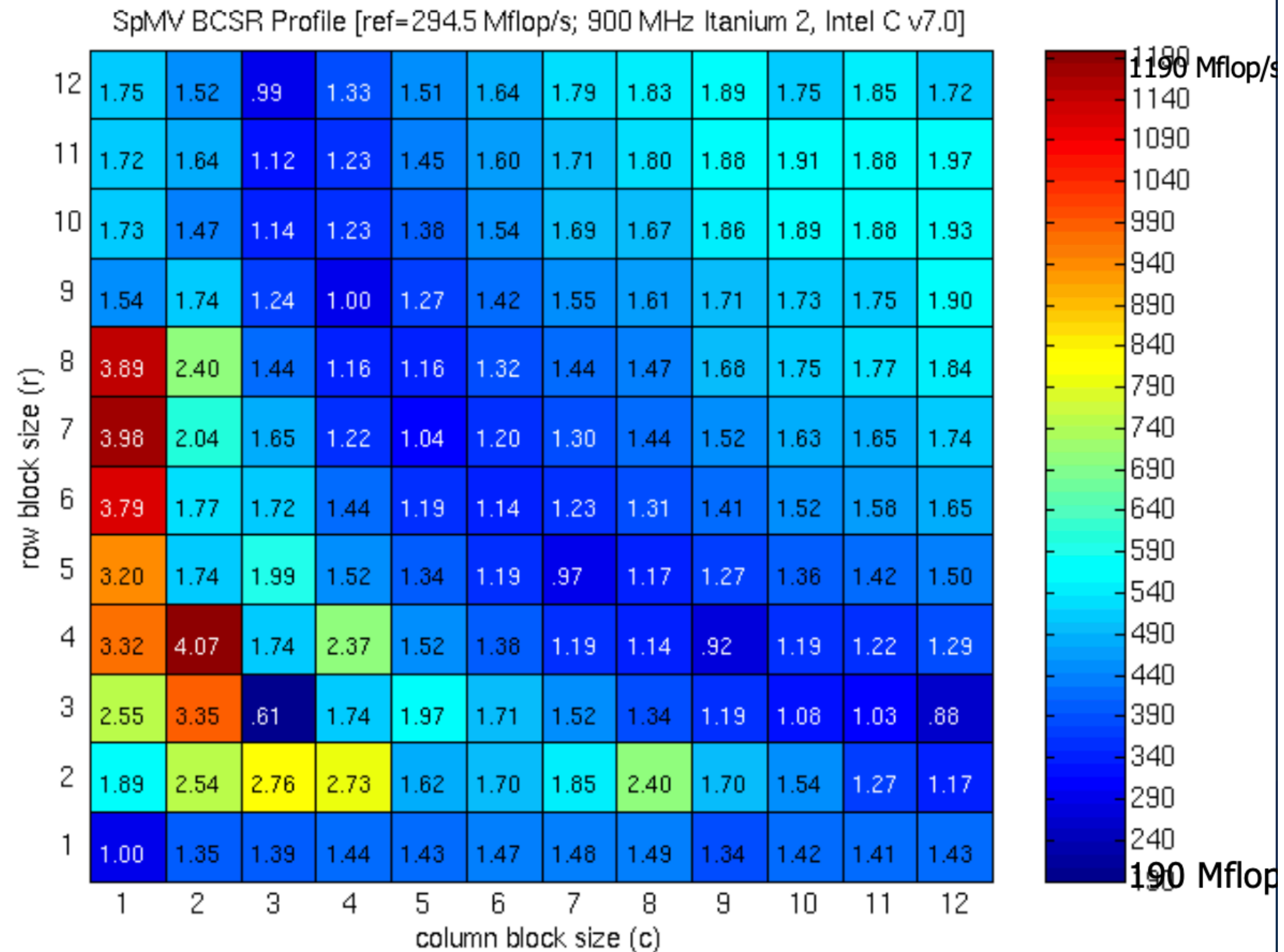**over all blockings**

Vuduc et al. show that when the fill was known exactly, performance of the resulting blocking scheme was **optimal or near optimal** (within 5%) [V04]
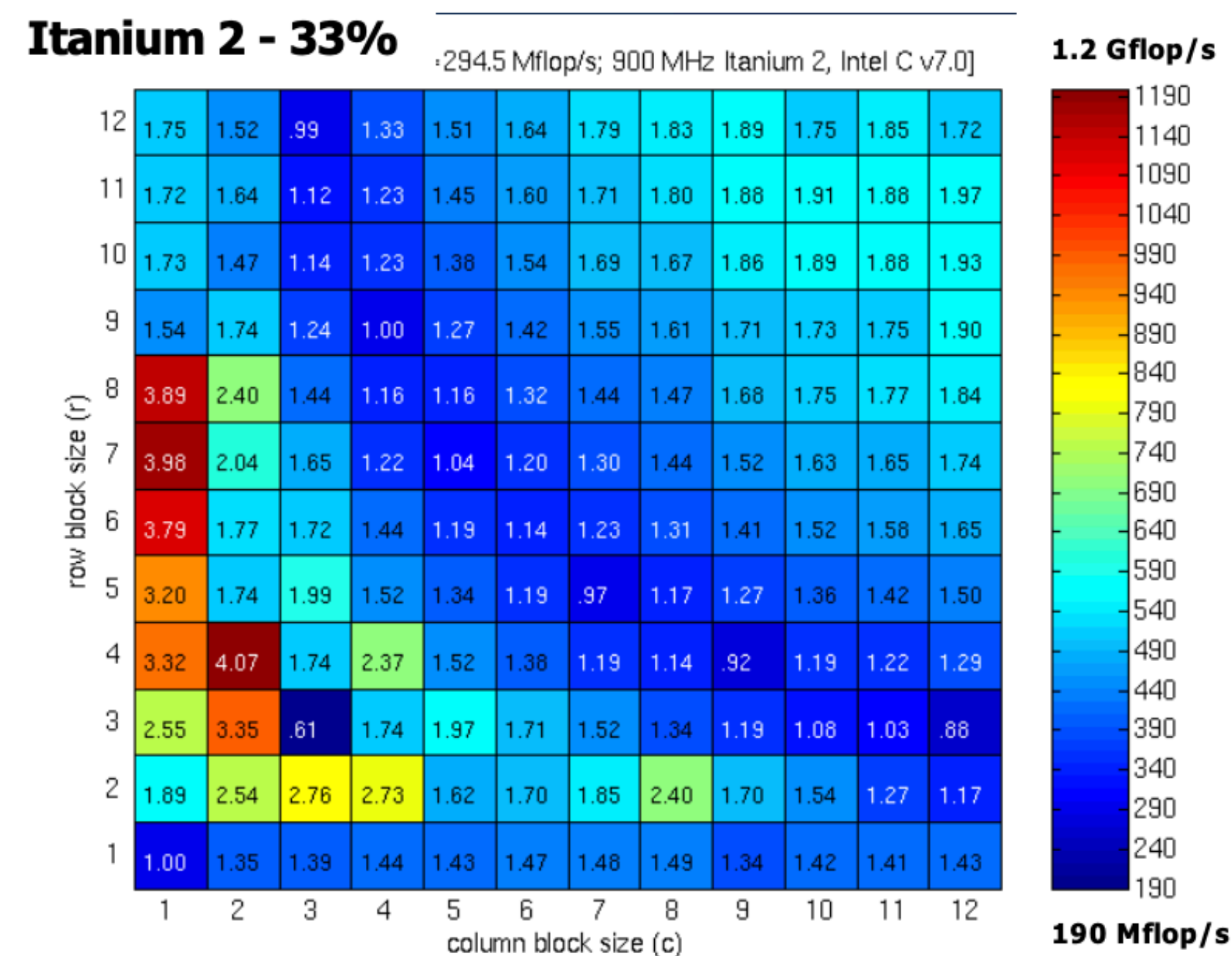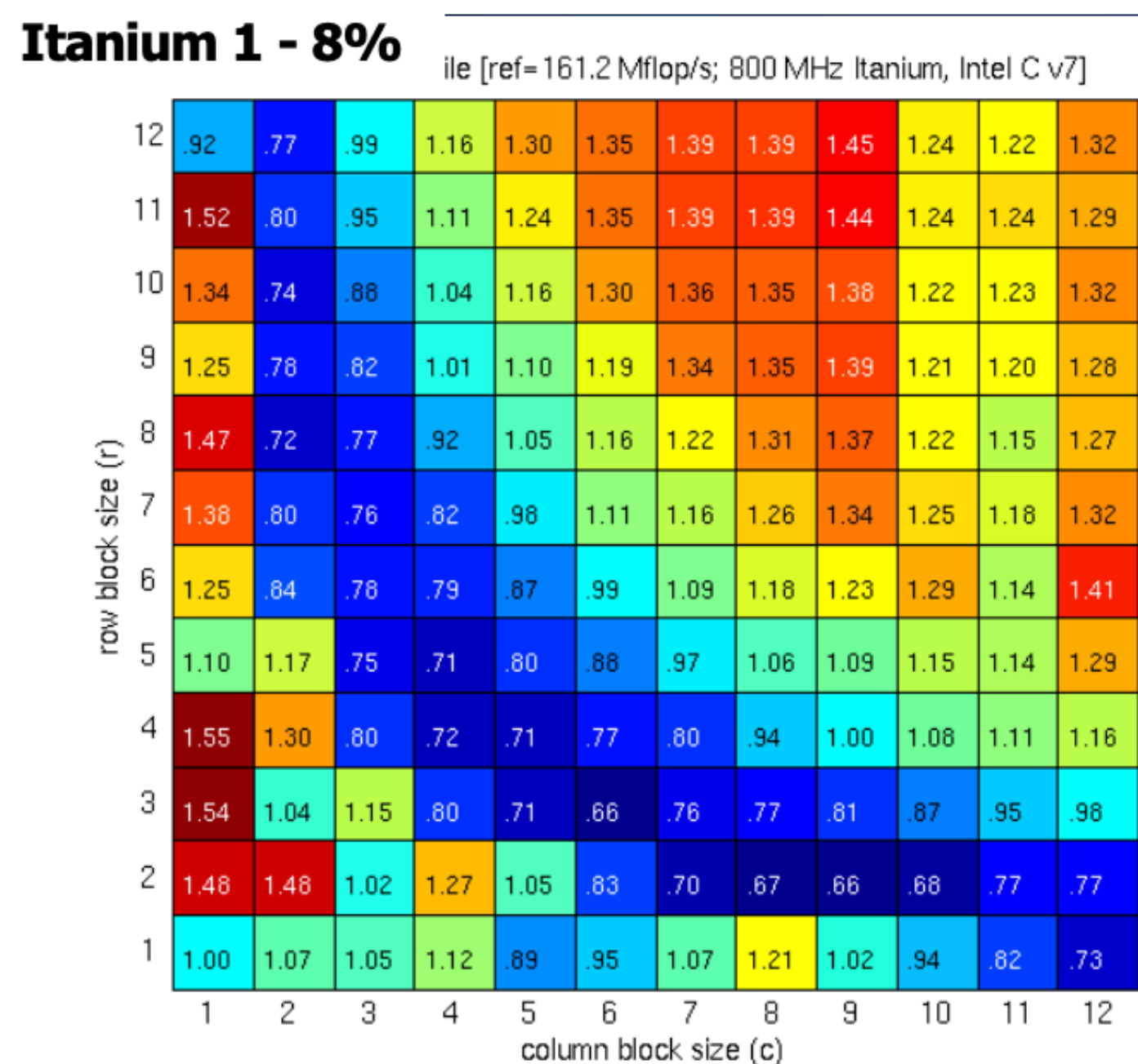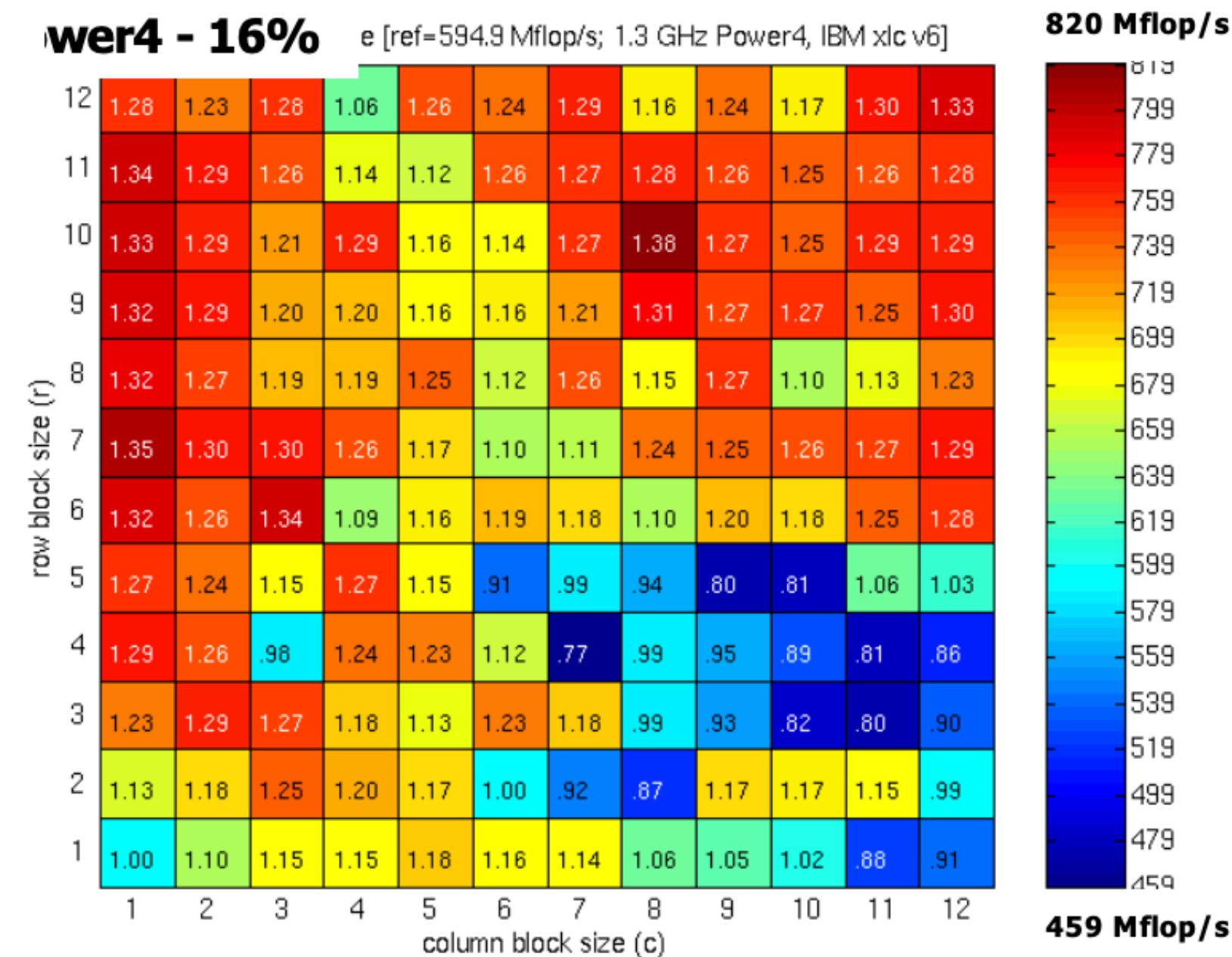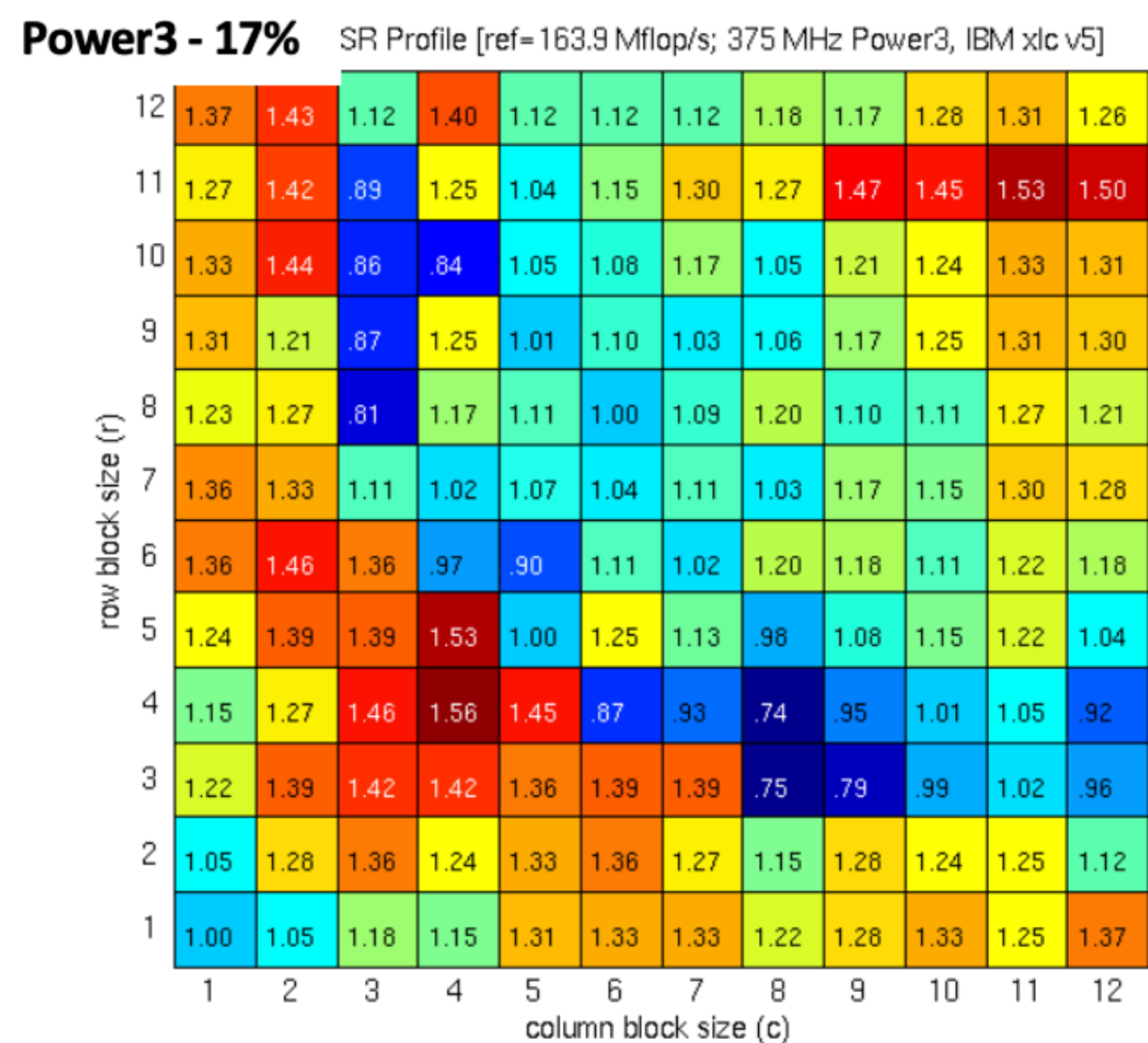
"Optimizing the Performance of Sparse Matrix-Vector Multiplication," Im, 2000.
"Automatic Performance Tuning of Sparse Matrix Kernels," Vuduc, 2003.
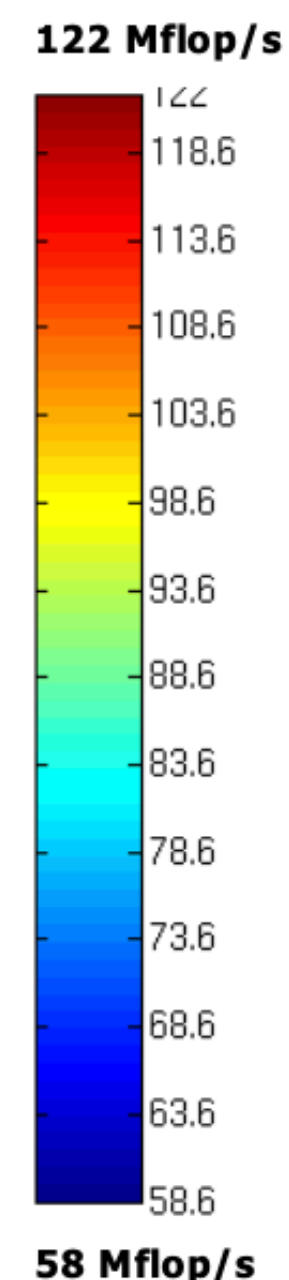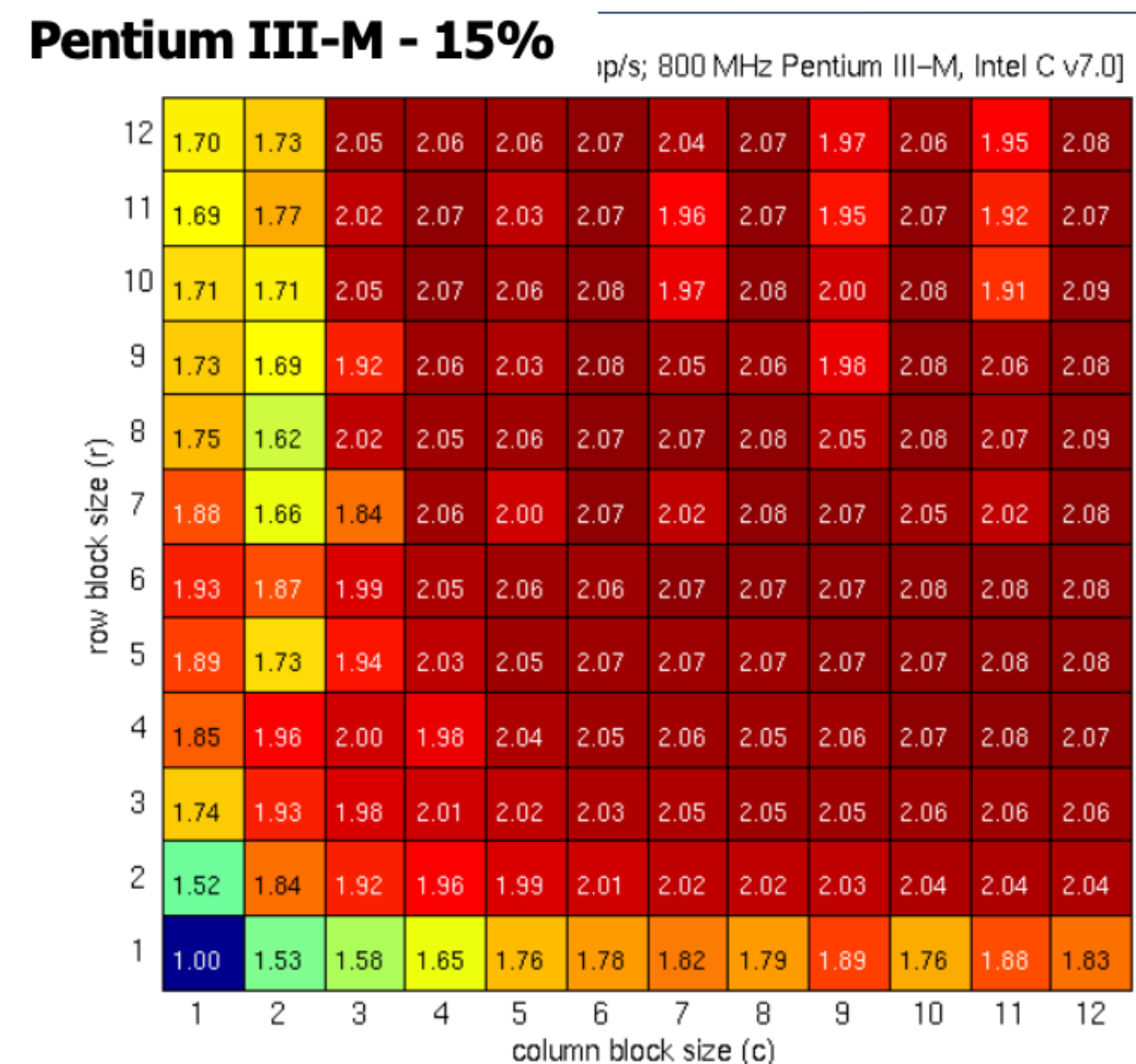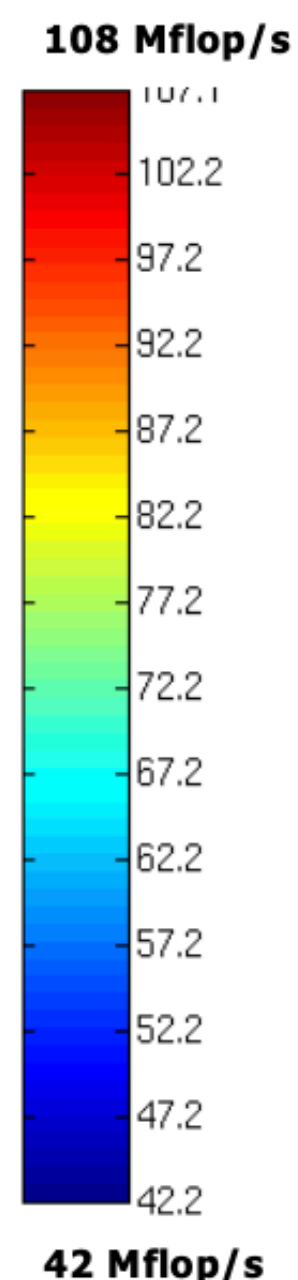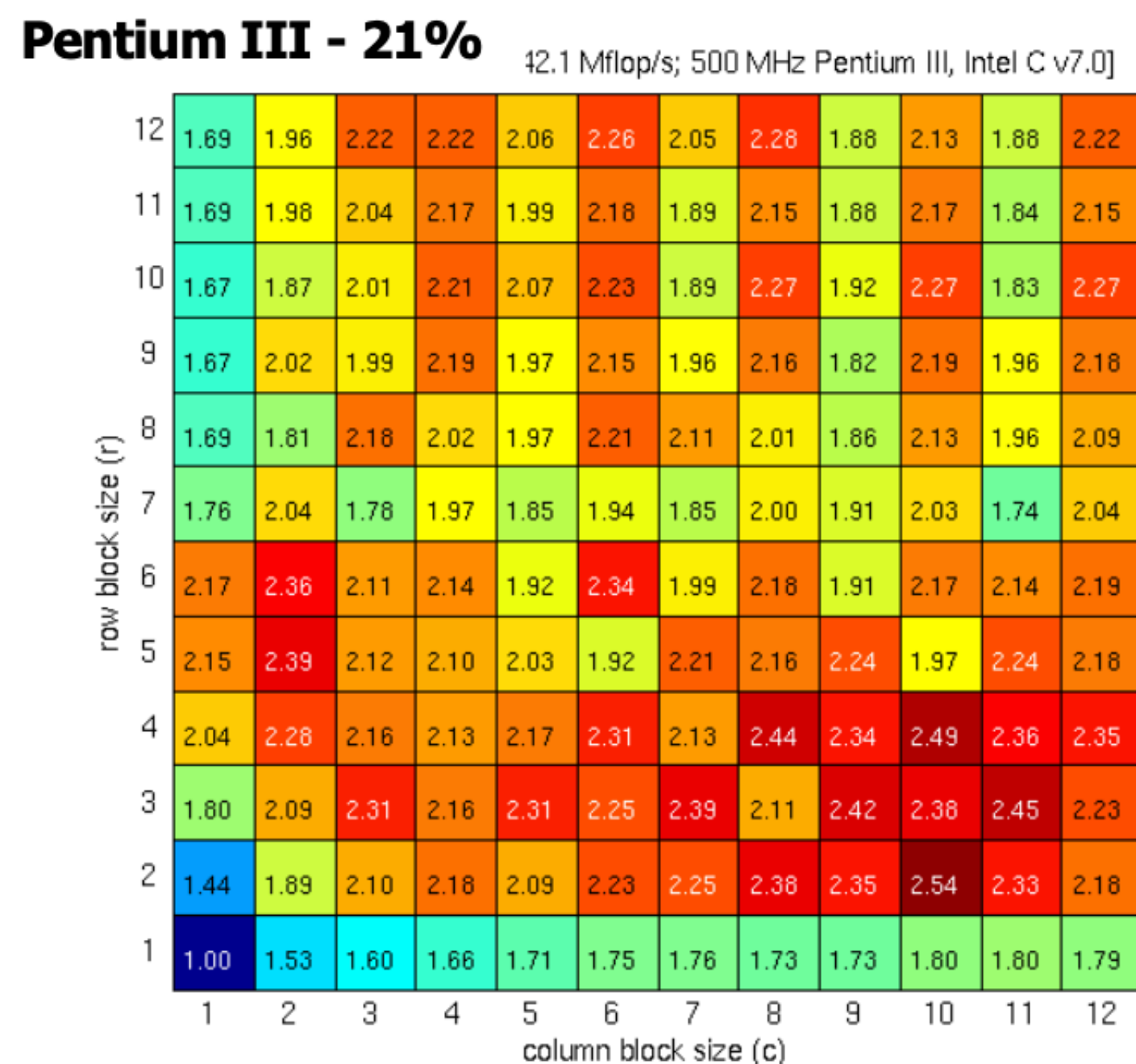
# Register Profile: dense matrix in sparse format



SpMV BCSR Profile [ref=294.5 Mflop/s; 900 MHz Itanium 2, Intel C v7.0]

| row block size (r) | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 12 | 1.75 | 1.52 | .99 | 1.33 | 1.51 | 1.64 | 1.79 | 1.83 | 1.89 | 1.75 | 1.85 | 1.72 |
| 11 | 1.72 | 1.64 | 1.12 | 1.23 | 1.45 | 1.60 | 1.71 | 1.80 | 1.88 | 1.91 | 1.88 | 1.97 |
| 10 | 1.73 | 1.47 | 1.14 | 1.23 | 1.38 | 1.54 | 1.69 | 1.67 | 1.86 | 1.89 | 1.88 | 1.93 |
| 9 | 1.54 | 1.74 | 1.24 | 1.00 | 1.27 | 1.42 | 1.55 | 1.61 | 1.71 | 1.73 | 1.75 | 1.90 |
| 8 | 3.89 | 2.40 | 1.44 | 1.16 | 1.16 | 1.32 | 1.44 | 1.47 | 1.68 | 1.75 | 1.77 | 1.84 |
| 7 | 3.98 | 2.04 | 1.65 | 1.22 | 1.04 | 1.20 | 1.30 | 1.44 | 1.52 | 1.63 | 1.65 | 1.74 |
| 6 | 3.79 | 1.77 | 1.72 | 1.44 | 1.19 | 1.14 | 1.23 | 1.31 | 1.41 | 1.52 | 1.58 | 1.65 |
| 5 | 3.20 | 1.74 | 1.99 | 1.52 | 1.34 | 1.19 | .97 | 1.17 | 1.27 | 1.36 | 1.42 | 1.50 |
| 4 | 3.32 | 4.07 | 1.74 | 2.37 | 1.52 | 1.38 | 1.19 | 1.14 | .92 | 1.19 | 1.22 | 1.29 |
| 3 | 2.55 | 3.35 | .61 | 1.74 | 1.97 | 1.71 | 1.52 | 1.34 | 1.19 | 1.08 | 1.03 | .88 |
| 2 | 1.89 | 2.54 | 2.76 | 2.73 | 1.62 | 1.70 | 1.85 | 2.40 | 1.70 | 1.54 | 1.27 | 1.17 |
| 1 | 1.00 | 1.35 | 1.39 | 1.44 | 1.43 | 1.47 | 1.48 | 1.49 | 1.34 | 1.42 | 1.41 | 1.43 |

column block size (c)

Color scale: 1190 Mflop/s, 1140, 1090, 1040, 990, 940, 890, 840, 790, 740, 690, 640, 590, 540, 490, 440, 390, 340, 290, 240, 190 Mflop/s

# Register Profiles



Power3 - 17%   SR Profile [ref=163.9 Mflop/s; 375 MHz Power3, IBM xlc v5]

Power4 - 16%   [ref=594.9 Mflop/s; 1.3 GHz Power4, IBM xlc v6]

Itanium 1 - 8%   [ref=161.2 Mflop/s; 800 MHz Itanium, Intel C v7]

Itanium 2 - 33%   [ref=294.5 Mflop/s; 900 MHz Itanium 2, Intel C v7.0]

# Register Profiles

# Heuristic Tuning vs Best



Accuracy of the Tuning Heuristics [Itanium 2]

Legend:
- ● Best
- □ Heuristic v2
- ◄ Exact fill
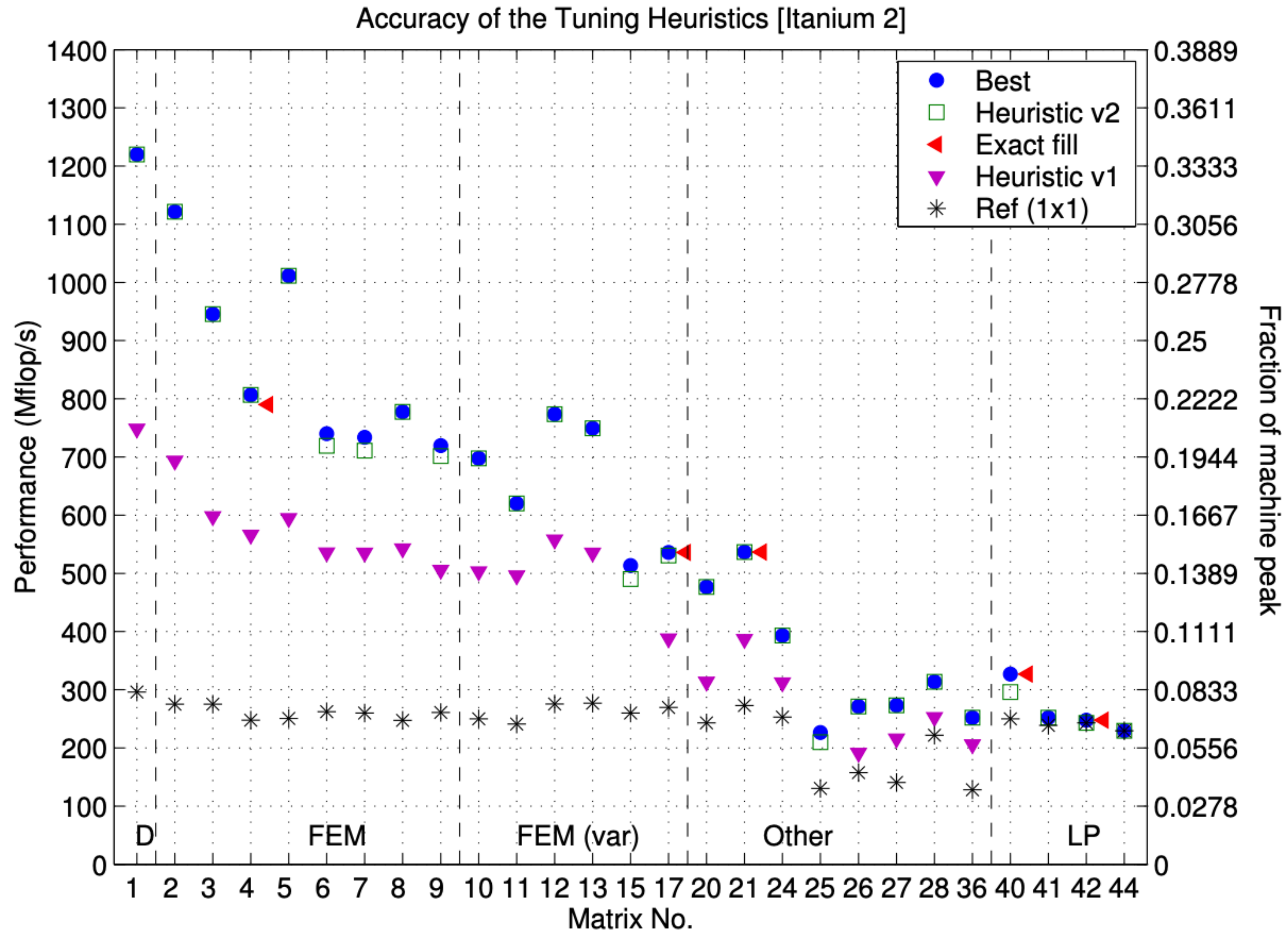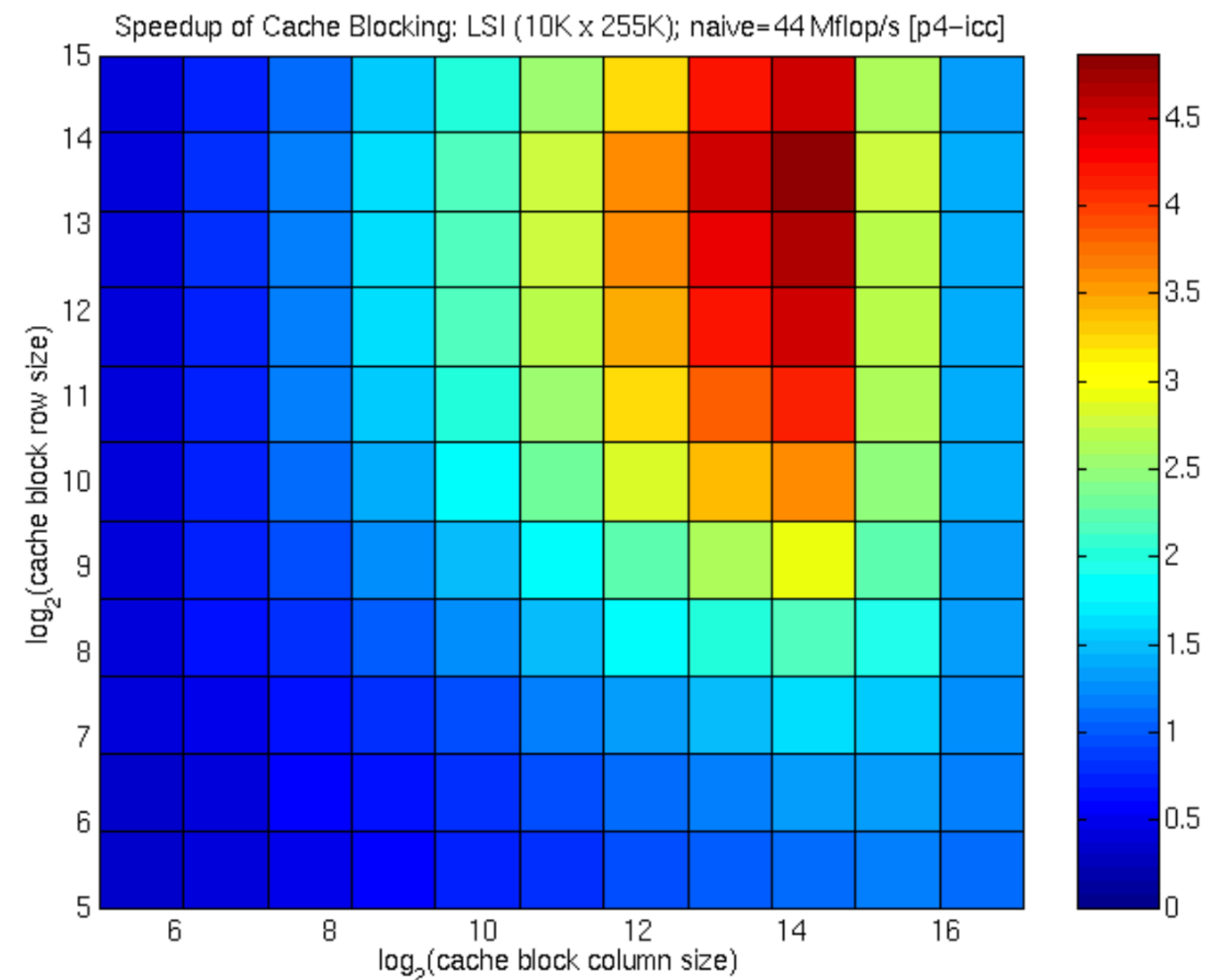- ▼ Heuristic v1
- ✳ Ref (1x1)

# What about cache blocking?

- For CSR, dot-product, re-use opportunity is only in the x vector
- Matrices that are have some locality and "well ordered" e.g., near diagonal have good re-use
- Cache blocking can help for "short wide" matrices both on serial and SMPs



Speedup of Cache Blocking: LSI (10K x 255K); naive = 44 Mflop/s [p4–icc]

Matrix A
- 100k x 255k
- 3.7M non-zeros

Baseline: 44 Mflop/s

Best block size & performance:
- 16k x 16k
- 210 Mflop/s

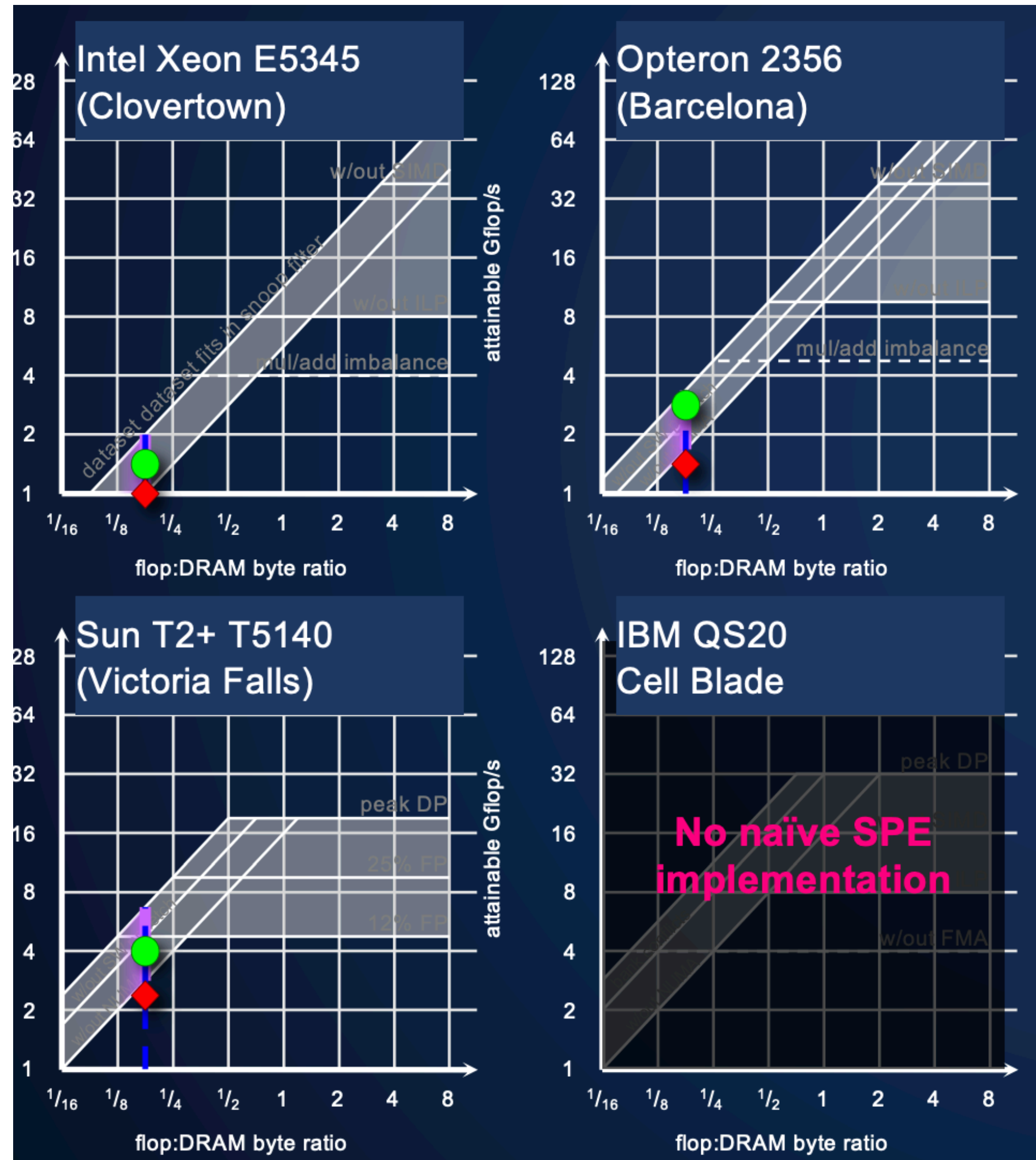# Roofline model for SpMV
# (out-of-the-box parallel)



Two unit stride streams

Inherent FMA

No instruction-level or data-level parallelism

Naïve compulsory: flop:byte < 0.166
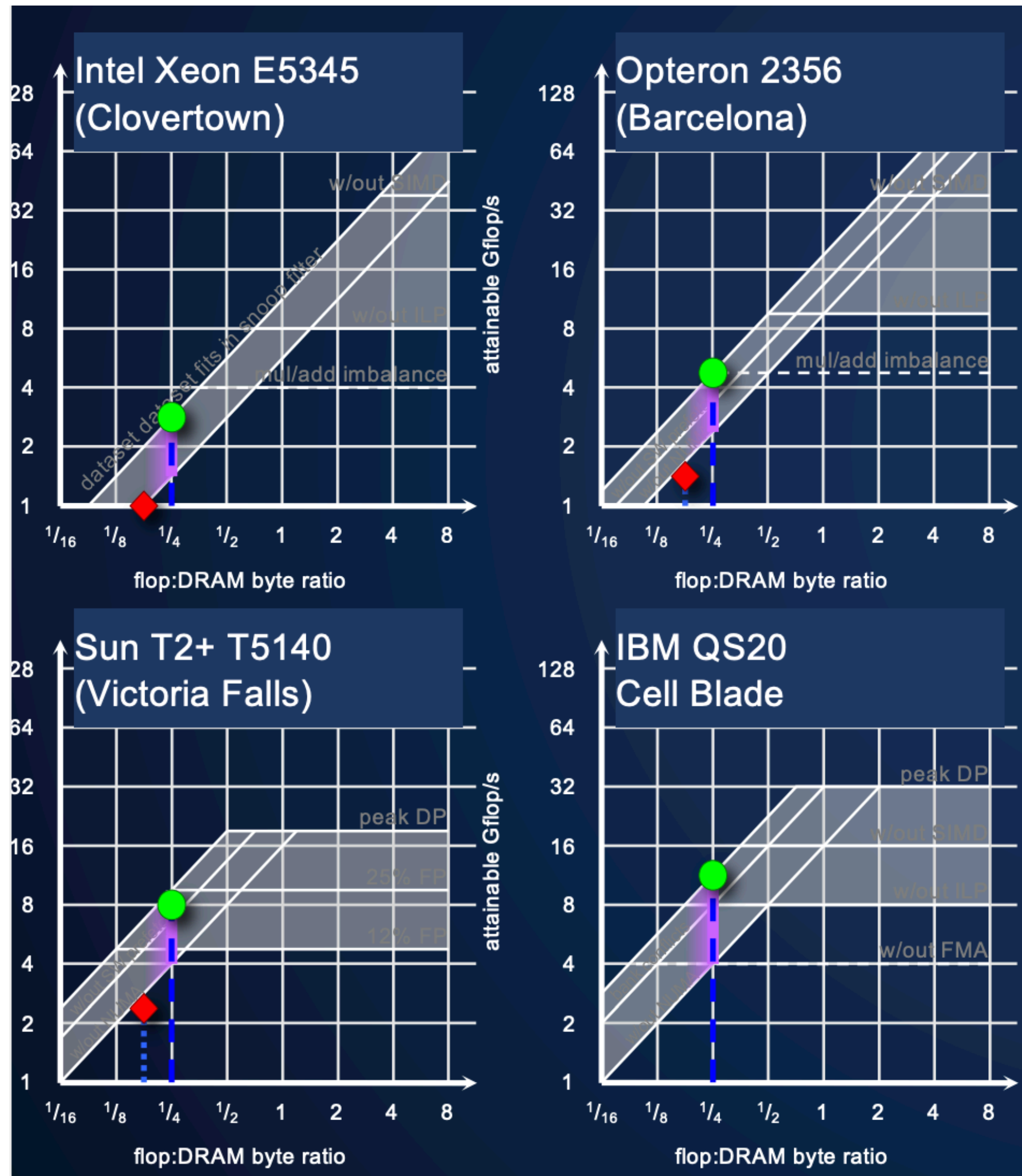
# Roofline model for SpMV
# (NUMA & SW prefetch)



Compulsory flop:byte ~ 0.166

Utilize all memory channels

Add software prefetching

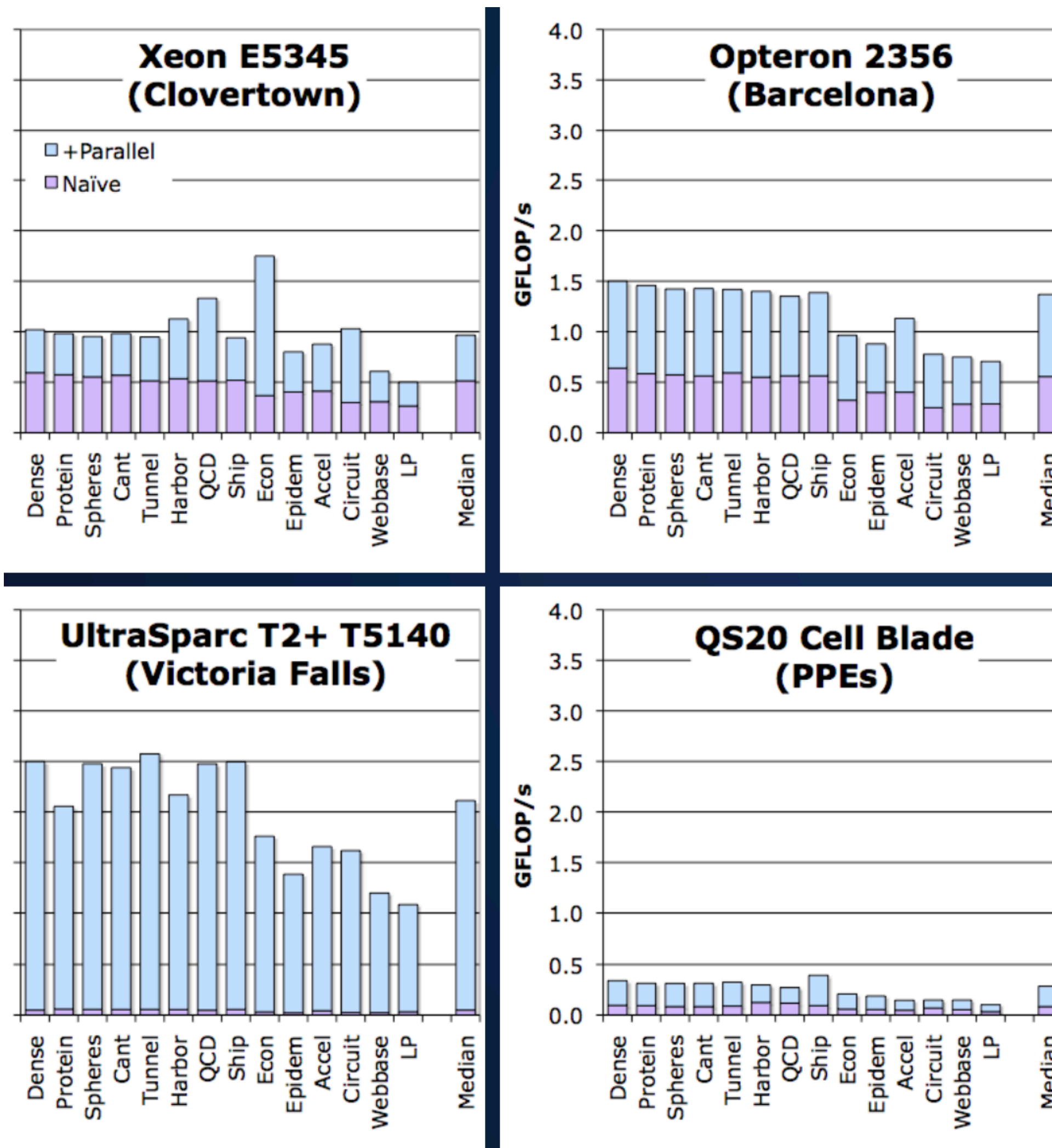# Roofline model for SpMV
# (matrix compression)



Inherent FMA

Register blocking improves ILP, DLP, flop:byte ratio

# SpMV Performance
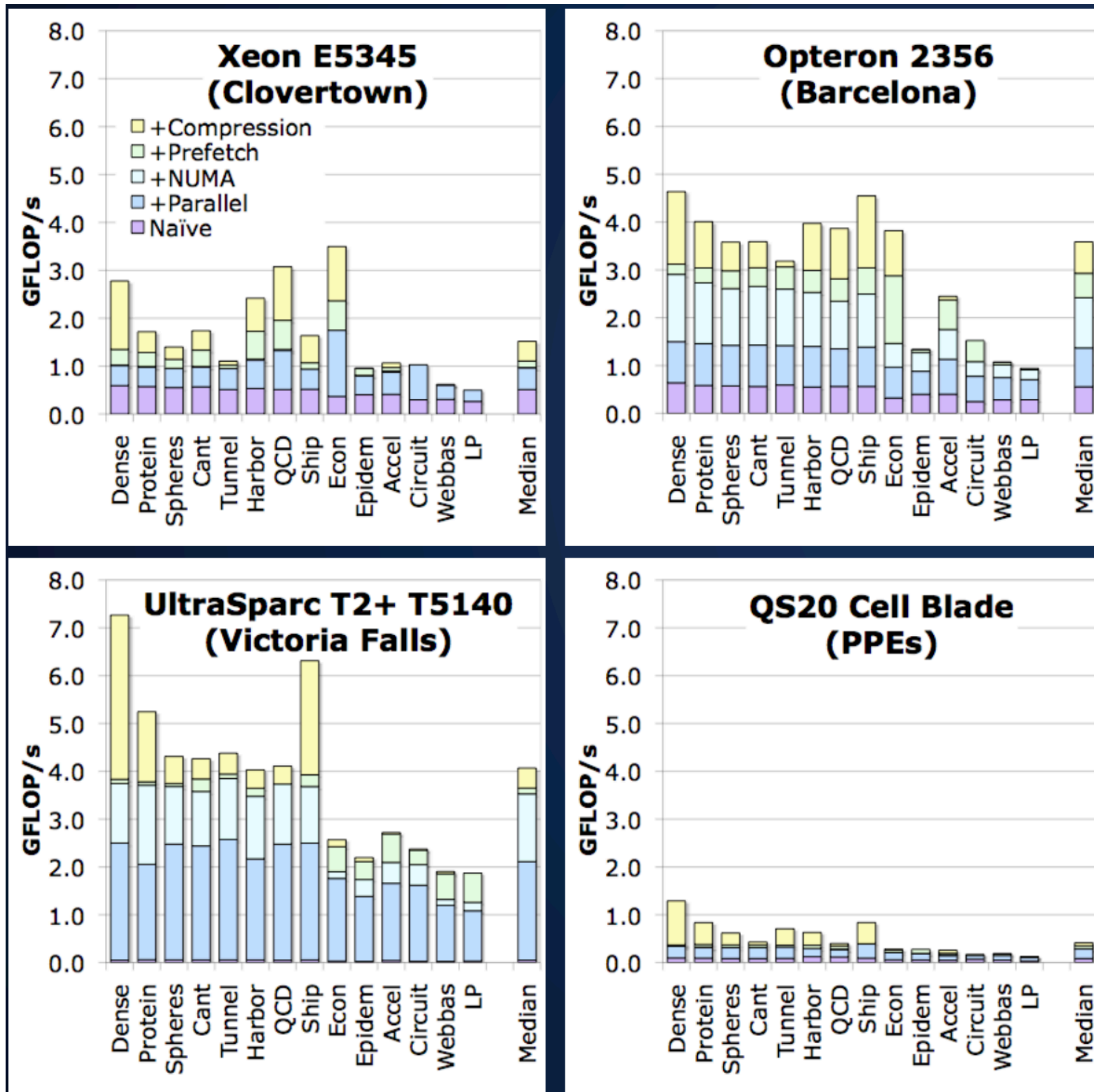# (simple parallelization)



Out-of-the box SpMV performance on a suite of 14 matrices

Simplest solution = parallelization by rows

Scalability isn't great - Can we do better?

# SpMV Performance
# (matrix compression)



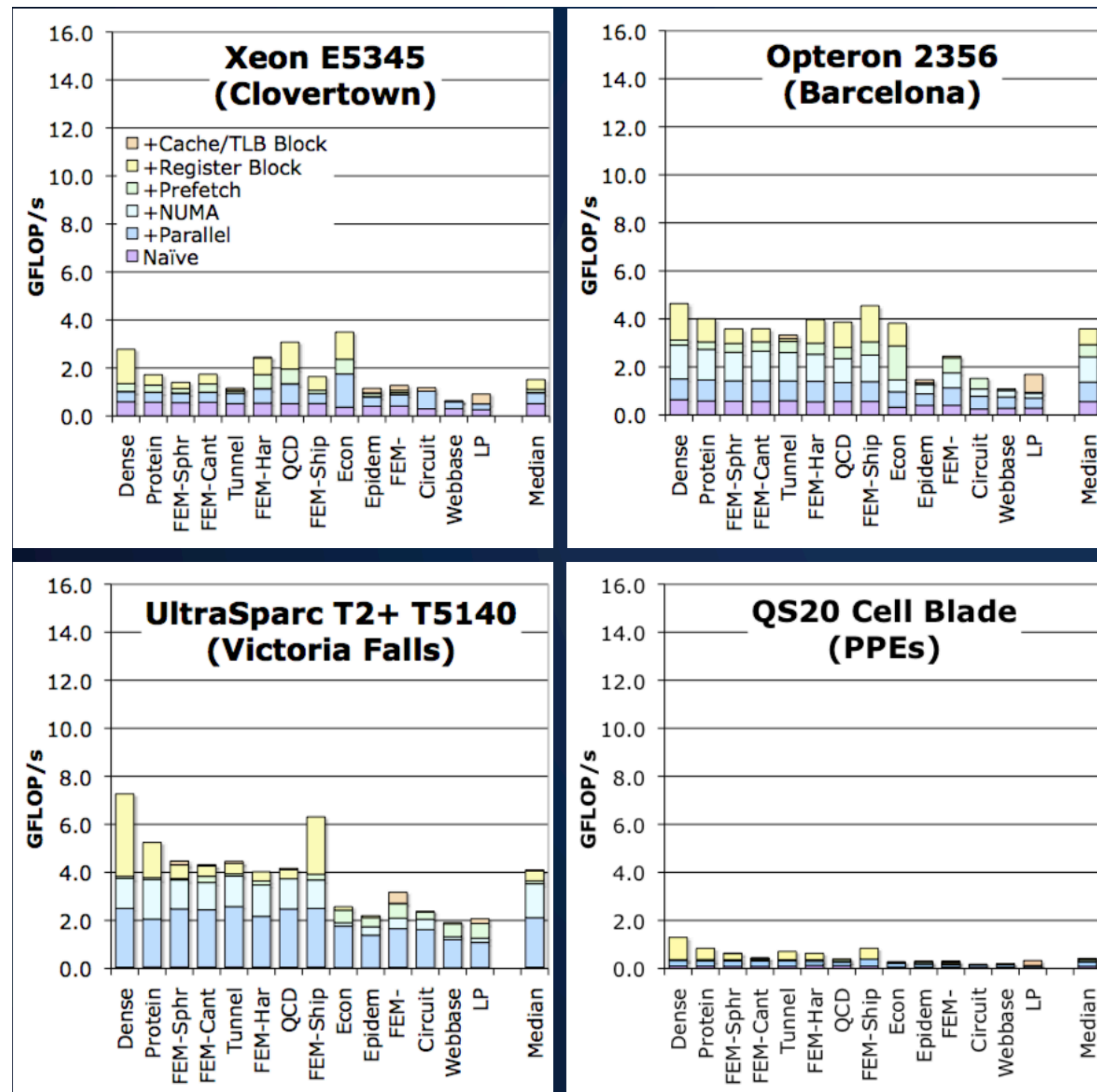After maximizing memory bandwidth, the only hope is to minimize memory traffic.

Compression: exploit
- register blocking
- other formats
- smaller indices

Benefit is matrix-dependent.

Register blocking enables efficient software prefetching (one per cache line)

# SpMV Performance
# (cache and TLB blocking)



Fully auto-tuned SpMV performance across the suite of matrices

Why do some optimizations work better on some architectures?
- Matrices with naturally small working sets
- Architectures with giant caches

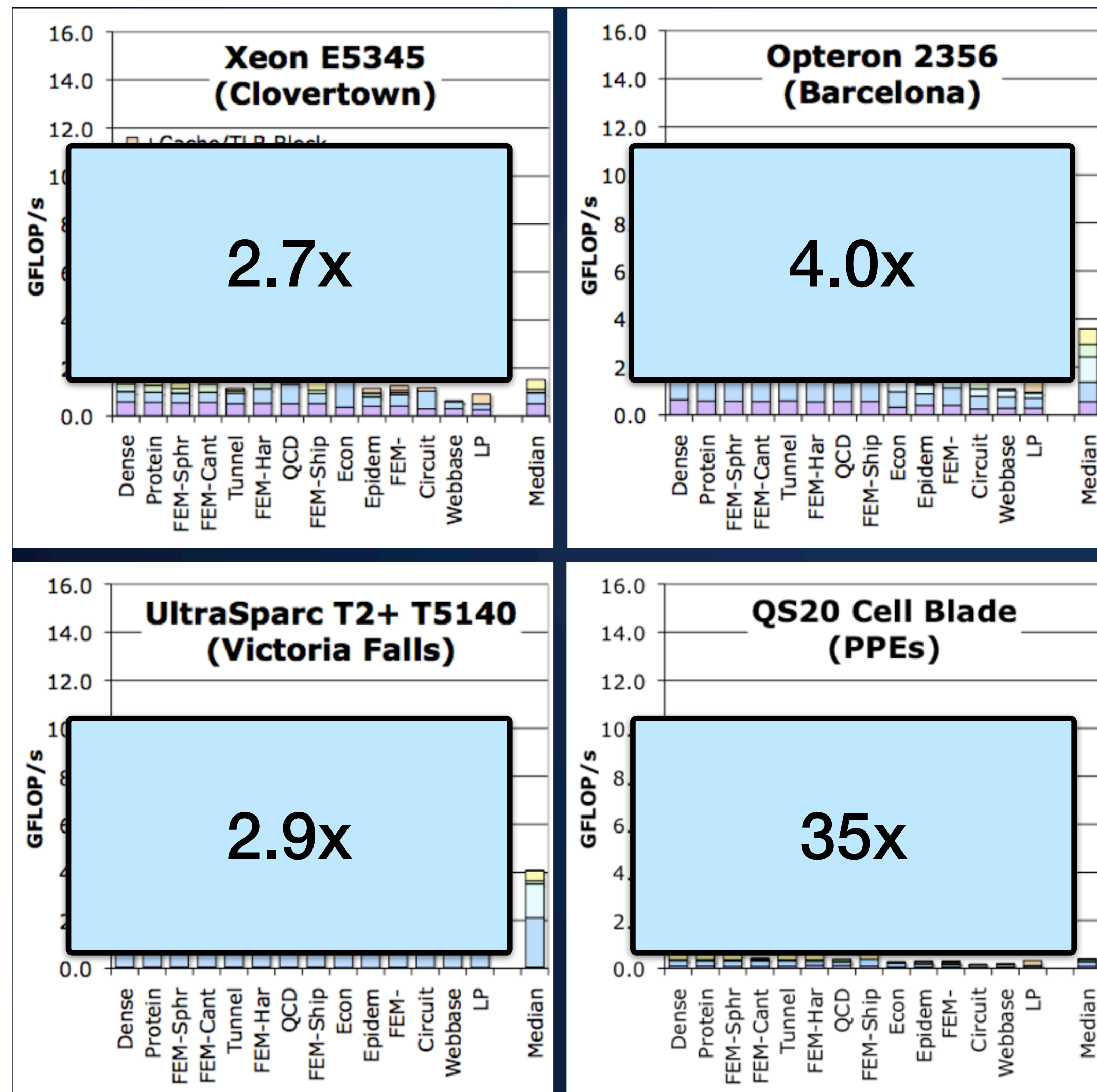# SpMV Performance
# (cache and TLB blocking)



Fully auto-tuned SpMV performance across the suite of matrices

Why do some optimizations work better on some architectures?
- Matrices with naturally small working sets
- Architectures with giant caches

# Summary

- Tuning for sparse matrices: harder than dense ones
- SpMV: benefits lower due to **low Computational Intensity** (read the matrix)
- **Register blocking** and other "compression" can be a big win
- Cache blocking less so; other low level tuning (e.g., prefetch) some
- For distributed memory, **reordering** (e.g., graph partitioning) important
- **Autotuning** possible, but depends on sparsity structure; hybrid offline / online tuning
- After tuning SpMV should be **memory bandwidth limited**
- Optimizing at **higher-level algorithms** (across iterations) can improve reuse.