





Performance Engineering for Sparse Linear Solvers Half-Day Tutorial at SC25

Christie L. Alappat, Erlangen National High Performance Computing Center

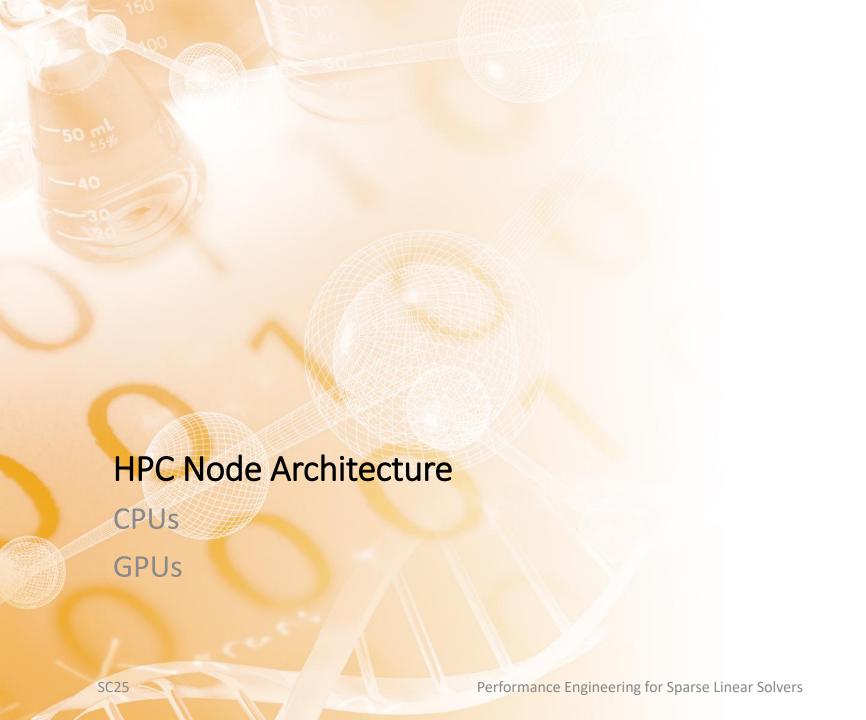
Jonas Thies, TU Delft

Georg Hager, Erlangen National High Performance Computing Center

Hartwig Anzt, TU München

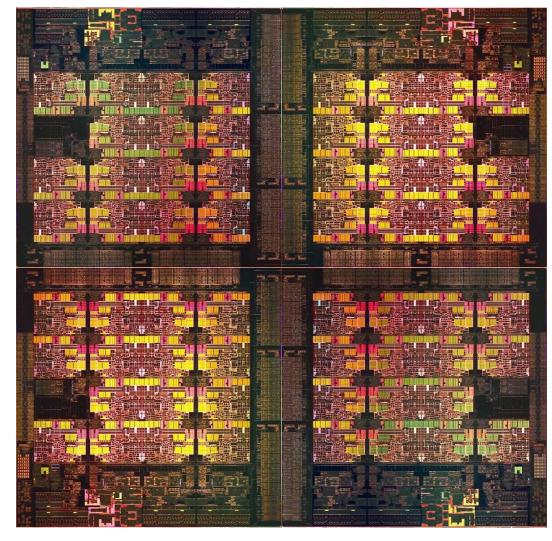
Tutorial Agenda

- Brief introduction to node-level computer architecture
- Performance modeling with the Roofline model
- Sparse matrix-vector multiplication (SpMV) performance, sparsematrix data formats, and Roofline modeling of SpMV
- The Conjugate Gradient (CG) algorithm
- Preconditioning and preconditioned CG (PCG)
- Accelerating matrix power kernels (MPK) by cache blocking
- Optional: distributed-memory SpMV and MPK cache blocking



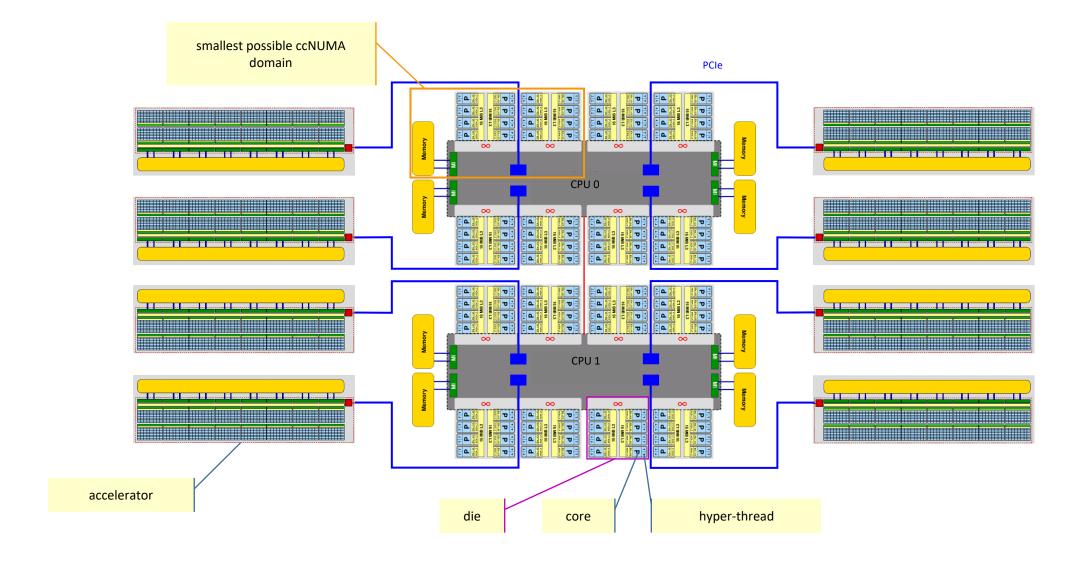
Multi-core today: Intel Xeon Sapphire Rapids (2023)

- Xeon "Sapphire Rapids" (Platinum/Gold/Silver/Bronze):
 Up to 60 cores running at 1.7+ GHz
 (+ "Turbo Mode" 4.8 GHz),
- "Intel 7" process / up to 350 W
- Multi-die package (4 chips)
- Clock frequency: flexible ©



https://www.techpowerup.com/292204/intel-sapphire-rapids-xeon-4-tile-mcm-annotated

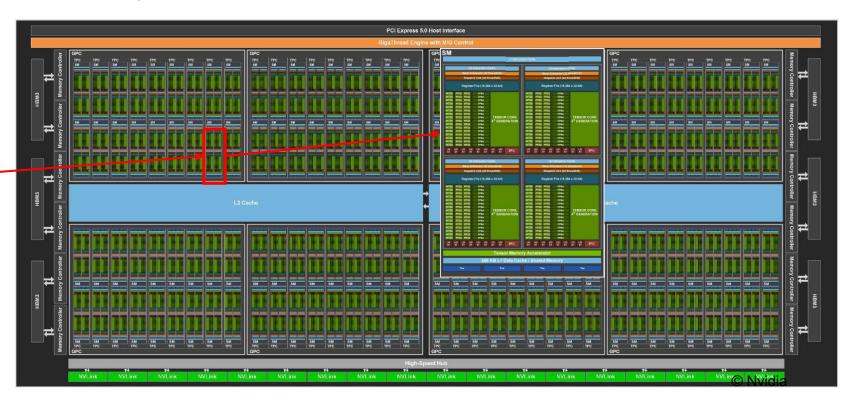
A more current example – lots of "topology"!



Nvidia H100 "Hopper" SXM5 specs

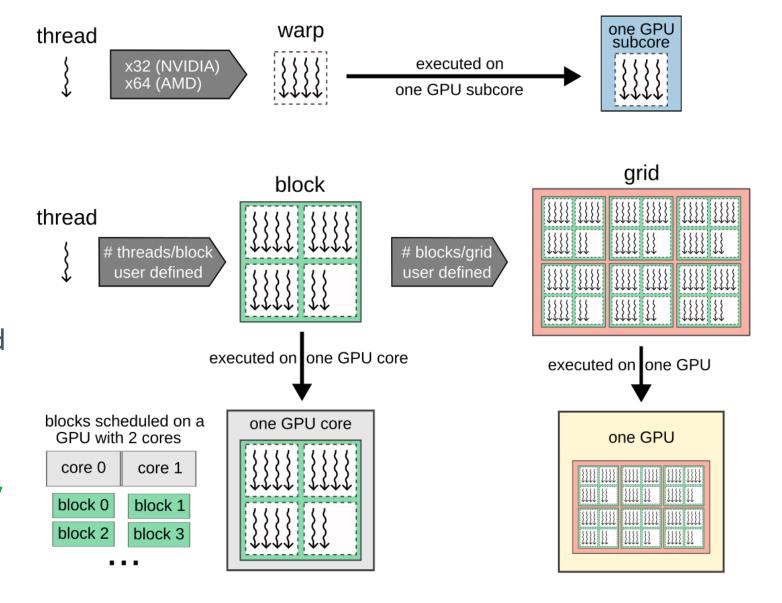
Architecture

- 80 B Transistors
- ~ 1.8 GHz clock speed
- ~ 144 "SM" units
 - 128 SP "cores" each (FMA)
 - 64 DP "cores" each (FMA)
 - 4 "Tensor Cores" each
 - 2:1 SP:DP performance
- ~ 34 TFlop/s DP peak (FP64)
- 50 MiB L2 Cache
- 80 GB HBM3
- MemBW ~ 3300 GB/s (theoretical)
- MemBW ~ 3000 GB/s (measured)



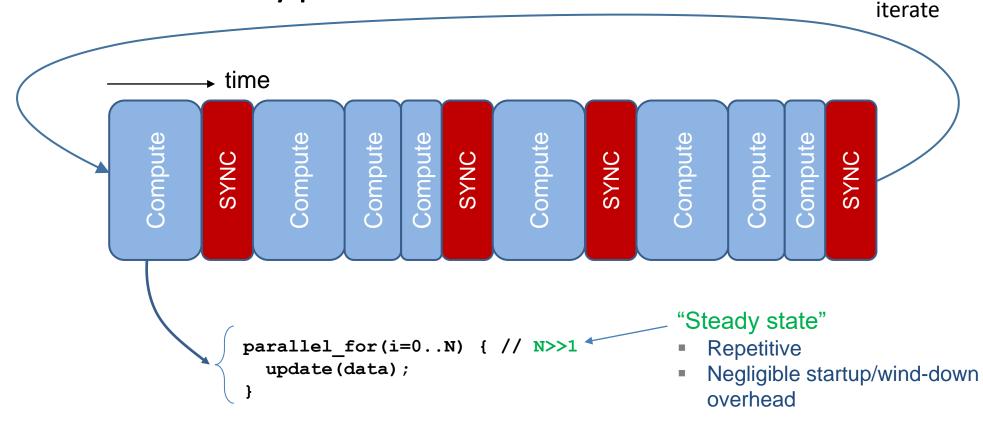
Threads, warps, grids on NVIDIA GPUs

- Threads are organized in warps (32 threads)
- Warps make up a block
- A block executes on one "core" (SM)
- Scheduling happens automatically
- Developer must map thread IDs to "work items" (iterations)
- Threads in a warp should access consecutive memory addresses



The Roofline Model Bottleneck-based thinking Simple models for single loops Multiple loops

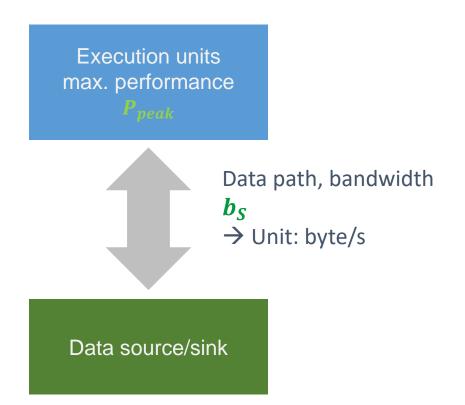
Structure of typical solver code



Runtime model: T = f(\$stuff, \$hardware)

A simple two-bottleneck model of loop code execution

Simplistic view of the hardware:



Simplistic view of the software:

```
! may be multiple levels
do i = 1,<sufficient>
        <complicated stuff doing
        N flops causing
        V bytes of data transfer>
enddo
```

Computational intensity $I = \frac{N}{V}$ \rightarrow Unit: flop/byte

Which takes longer?

- Data transfer
- Work execution

Predicting the (minimum) runtime of a loop

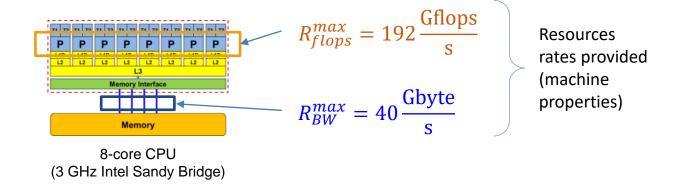
Two bottlenecks:

```
#pragma omp parallel for
for(i=0; i<10<sup>7</sup>; ++i)
    a[i] = a[i] + s * c[i];
```

Resources needed (code properties)

$$W_{flops} = 2 \times 10^7 \text{ flops}$$

 $W_{BW} = 3 \times 8 \times 10^7 \text{ bytes}$



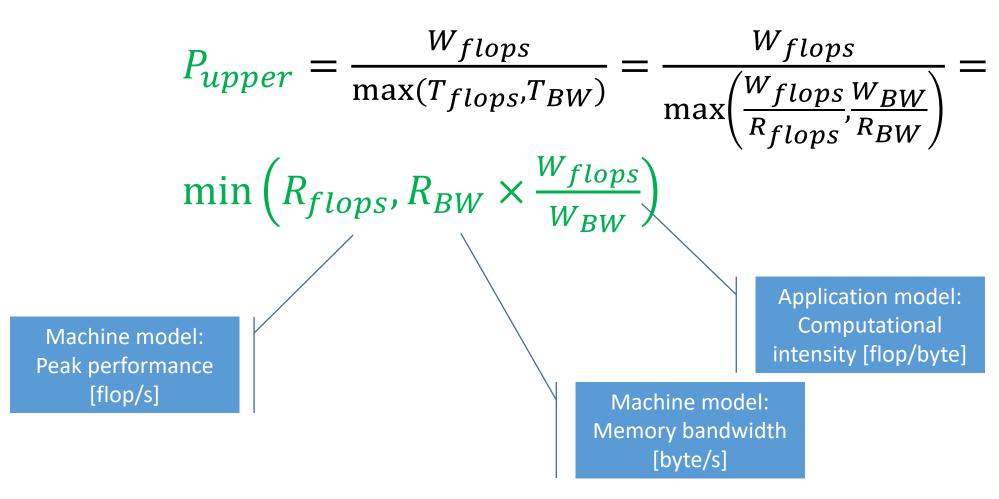
Full-overlap assumption:

$$T_{flops} = \frac{2 \times 10^7 \text{ flops}}{192 \frac{\text{Gflops}}{\text{s}}} = 104 \,\mu\text{s}$$
 $T_{BW} = \frac{2.4 \times 10^8 \text{ bytes}}{40 \frac{\text{Gbyte}}{\text{s}}} = 6.0 \text{ ms}$

$$T_{\min} = \max(T_{flops}, T_{BW})$$

= 6 ms

From time to performance



 $f_{1/2}$: A parameter to characterize memory and communication bottlenecks. Parallel Computing 10, 277-286 (1989). DOI: 10.1016/0167-8191(89)90100-2

"Roofline"!?

S. Williams: <u>Auto-tuning Performance on Multicore Computers</u>. UCB Technical Report No. UCB/EECS-2008-164. PhD thesis (2008)

Common nomenclature:

 R_{flops}

 $\rightarrow P_{peak}$ peak performance

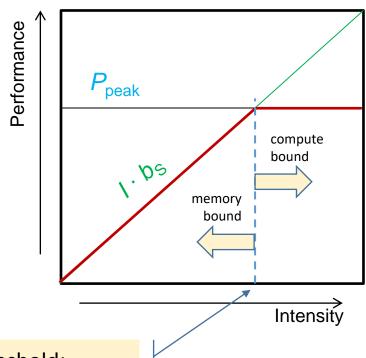
 R_{BW}

 $\rightarrow b_{\rm S}$ memory bandwidth

 $\frac{W_{flops}}{W_{BW}}$

 \rightarrow I computational intensity





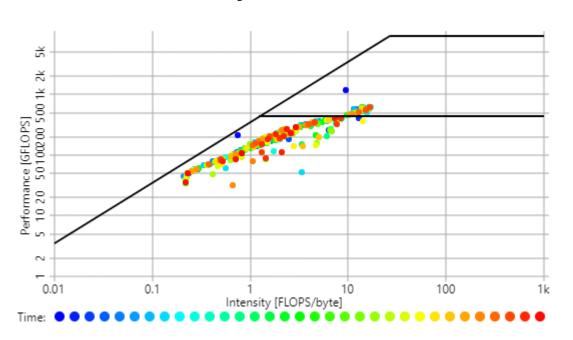
Threshold: ≈ 10-15 F/B for current Server CPUs/GPUs

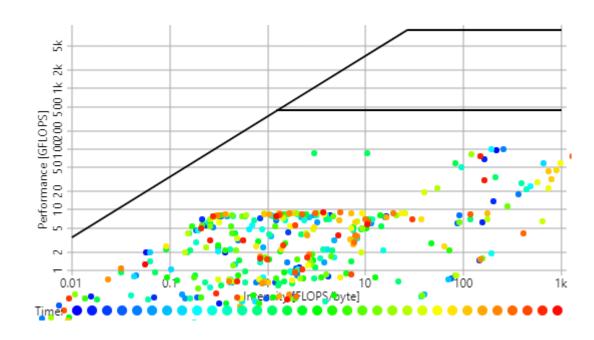
Hands-On: Exploring node topology and bandwidth

Diagnostic modeling

P_{max} P_{meas} N_{meas} V_{meas} Intensity

Two cluster jobs...





Roofline: Simple Examples Dense linear algebra Sparse linear algebra Simple solvers: CG SC25 Performance Engineering for Sparse Linear Solvers

Dense linear algebra

```
for(i=0; i<N; ++i)
    a[i] = a[i]+s*x[i];
    daxpy (BLAS-1)</pre>
```

```
for(i=0; i<N; ++i)
s += a[i]*b[i];</pre>
```

dot product (BLAS-1)

Roofline thinking:

What is the computational intensity?

```
for(k=0; k<NK;++k)
  for(l=0; l<NL; ++l)
   for(m=0; m<NM; ++m)
    y[k*NL+l] +=
        A[k*NM+m]*B[l*NM+m];
        dense MMM (BLAS-3)</pre>
```

```
for(r=0; r<NR; ++r)
  for(c=0; c<NC; ++c)
   y[r] += A[r*NC+c]*x[c];

  dense MVM (BLAS-2)
        dot-product style</pre>
```

Dense MVM

```
for(r=0; r<NR; ++r)
for(c=0; c<NC; ++c)
y[r] += A[r*NC+c]*x[c];
```

- One DP read from memory for each matrix entry
- x[] and y[] are read and updated from cache after 1st read
- > 8 byte and 2 flops per iteration

```
Computational intensity I = \frac{2 flop}{8 byte} = 0.25 \frac{flop}{byte}
```

Dense MMM?

 Blocking/unrolling techniques can increase intensity beyond the Roofline knee

→ peak performance achievable

Sparse Matrices and SpMV Sparse Matrix Formats Sparse Matrix Vector Product Parallelization

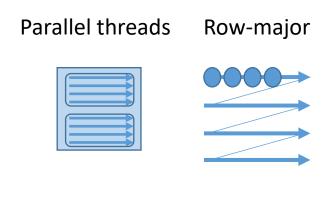
Input A, x, y Output $y = A \cdot x$

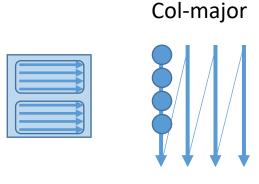
- Central building block in many complex algorithms:
 - Orthogonalization, power iteration in Page Rank, Power flow of a system ...
- Before we turn to sparse matrices, we recall how we store & handle dense matrices on parallel processors (i.e. GPUs)

```
Input A, x, y Output y = A \cdot x
```

```
__global__ void sgemv_rowmajor( ...)
{
    int row = blockIdx.x*blockDim.x + threadIdx.x;
    float sum = 0.0;
    if (row < n){
        for( int col=0; col<n; col++){
            sum += m[ row*n + col ] * x[ col ];
        }
        y[ row ] = alpha * sum;
    }
}
```

```
__global__ void sgemv_colmajor( ...)
{
    int row = blockldx.x*blockDim.x + threadIdx.x;
    float sum = 0.0;
    if (row < n){
        for( int col=0; col<n; col++){
            sum += m[ row + n*col ] * x[ col ];
        }
        y[ row ] = alpha * sum;
    }
}
```



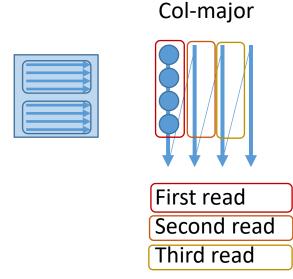


Input A, x, y Output $y = A \cdot x$

```
__global__ void sgemv_rowmajor( ...)
{
    int row = blockldx.x*blockDim.x + threadIdx.x;
    float sum = 0.0;
    if (row < n){
        for( int col=0; col<n; col++){
            sum += m[ row*n + col ] * x[ col ];
        }
        y[ row ] = alpha * sum;
    }
}
```

```
__global__ void sgemv_colmajor( ...)
{
    int row = blockldx.x*blockDim.x + threadIdx.x;
    float sum = 0.0;
    if (row < n){
        for( int col=0; col<n; col++){
            sum += m[ row + n*col ] * x[ col ];
        }
        y[ row ] = alpha * sum;
    }
}
```

Parallel threads Row-major



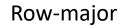
Input A, x, y Output $y = A \cdot x$

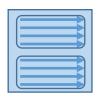
```
__global__ void sgemv_rowmajor( ...)
{
    int row = blockldx.x*blockDim.x + threadIdx.x;
    float sum = 0.0;
    if (row < n){
        for( int col=0; col<n; col++){
            sum += m[ row*n + col ] * x[ col ];
        }
        y[ row ] = alpha * sum;
    }
}
```

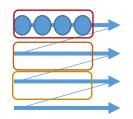
```
__global__ void sgemv_colmajor( ...)

{
    int row = blockldx.x*blockDim.x + threadIdx.x;
    float sum = 0.0;
    if (row < n){
        for( int col=0; col<n; col++){
            sum += m[ row + n*col ] * x[ col ];
        }
        y[ row ] = alpha * sum;
    }
}
```

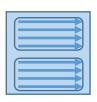
Parallel threads

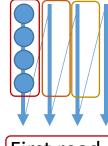


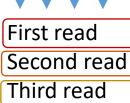


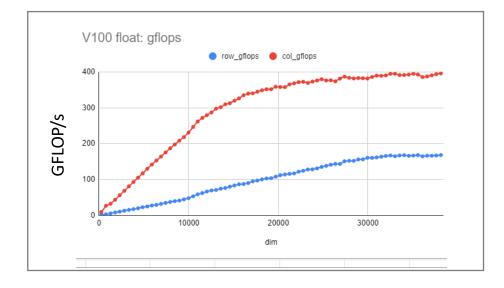


Col-major









...

Sparse Matrix Vector Multiplication

Input A, x, y Output $y = A \cdot x$

- Matrix A contains only few nonzero elements.
- Storing all entries results in large overhead (memory & computation).

Sparse Matrix Vector Multiplication

Input A, x, y Output $y = A \cdot x$

- Matrix A contains only few nonzero elements.
- Storing all entries results in large overhead (memory & computation).
- Idea: Store only nonzero elements [nz] explicitly.

$$A = \begin{pmatrix} 5.4 & 1.1 & 0 & 0 & 0 & 0 \\ 2.2 & 8.3 & 0 & 3.7 & 1.3 & 3.8 \\ 0 & 0 & 4.2 & 0 & 0 & 0 \\ 5.4 & 0 & 0 & 9.2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1.1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 8.1 \end{pmatrix}$$

 $value = [5.4 \ 1.1 \ 2.2 \ 8.3 \ 3.7 \ 1.3 \ 3.8 \ 4.2 \ 5.4 \ 9.2 \ 1.1 \ 8.1]$ Value

- Matrix A contains only few nonzero elements.
- Storing all entries results in large overhead (memory & computation).
- Idea: Store only nonzero elements [nz] explicitly.

Need to also store location of nonzero elements!

$$A = \begin{pmatrix} 5.4 & 1.1 & 0 & 0 & 0 & 0 \\ 2.2 & 8.3 & 0 & 3.7 & 1.3 & 3.8 \\ 0 & 0 & 4.2 & 0 & 0 & 0 \\ 5.4 & 0 & 0 & 9.2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1.1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 8.1 \end{pmatrix}$$

Memory footprint of COO format: nz(val) + 2*nz(int)

COO format

 $\operatorname{Input} A, x, y \quad \operatorname{Output} y = A \cdot x$

Matrix A continuous formula A Storing all ent Hands-on Exercise: Convert this matrix into COO format: Idea: Store or Need to $\left(\begin{array}{cccccccc}
0 & 0 & 0 & 4 & 2 & 0 \\
0 & 2 & 3 & 0 & 0 & 1 \\
0 & 0 & 0 & 0 & 0 & 0 \\
1 & 2 & 3 & 4 & 3 & 1 \\
0 & 2 & 0 & 0 & 0 & 1 \\
1 & 2 & 3 & 0 & 0 & 1
\end{array}\right)$ format: $value = \begin{bmatrix} 5 \end{bmatrix}$ $colidx = \begin{bmatrix} 0 \end{bmatrix}$ Compute the memory requirement (# vals + # int) $rowidx = \begin{bmatrix} 0 \end{bmatrix}$ **Row-index**

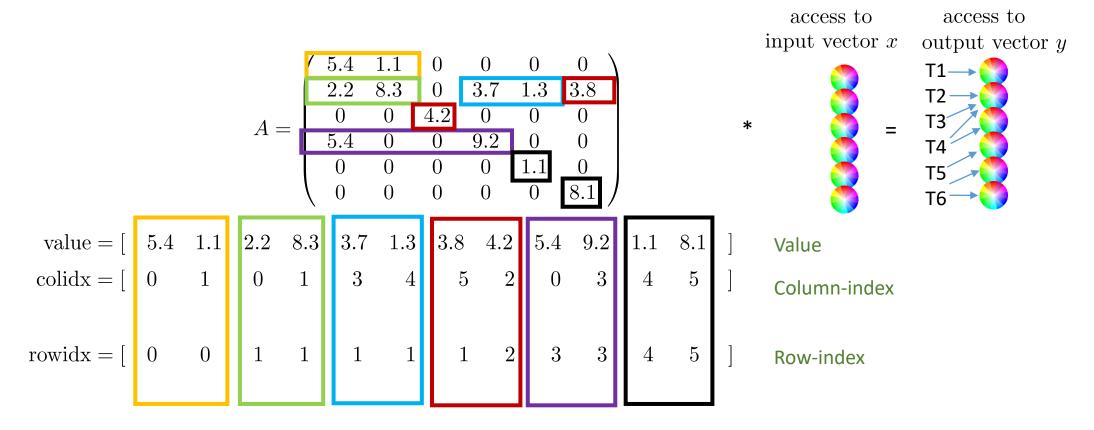
COO format

 $\operatorname{Input} A, x, y \quad \operatorname{Output} y = A \cdot x$

Matrix A continuous formula A Storing all ent Hands-on Exercise: Convert this matrix into COO format: Idea: Store or Need to $\left(\begin{array}{ccccccc}
0 & 0 & 0 & 4 & 2 & 0 \\
0 & 2 & 3 & 0 & 0 & 1 \\
0 & 0 & 0 & 0 & 0 & 0 \\
1 & 2 & 3 & 4 & 3 & 1 \\
0 & 2 & 0 & 0 & 0 & 1 \\
1 & 2 & 3 & 0 & 0 & 1
\end{array}\right)$ format: $value = \begin{bmatrix} 5 \end{bmatrix}$ $colidx = \begin{bmatrix} 0 \end{bmatrix}$ Compute the memory requirement (# vals + # int) $rowidx = \begin{bmatrix} 0 \end{bmatrix}$

Split nonzero elements into chunks and parallelize across chunks.

- Partial sums need synchronization / atomics to avoid write conflicts.
- Non-coalesced memory access (because row-major).



CSR (==CRS) format

 $\operatorname{Input} A, x, y \quad \operatorname{Output} y = A \cdot x$

- Matrix A contains only few nonzero elements.
- Storing all entries results in large overhead (memory & computation).
- Idea: Store only nonzero elements [nz] explicitly.

Need to also store location of nonzero elements!

$$A = \begin{pmatrix} 5.4 & 1.1 & 0 & 0 & 0 & 0 \\ 2.2 & 8.3 & 0 & 3.7 & 1.3 & 3.8 \\ 0 & 0 & 4.2 & 0 & 0 & 0 \\ 5.4 & 0 & 0 & 9.2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1.1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 8.1 \end{pmatrix}$$

Memory footprint of COO format: nz(val) + 2*nz(int)

CSR format

Input A, x, y Output $y = A \cdot x$

- Matrix A contains only few nonzero elements.
- Storing all entries results in large overhead (memory & computation).
- Idea: Store only nonzero elements [nz] explicitly.

Need to also store location of nonzero elements!

$$A = \begin{pmatrix} 5.4 & 1.1 & 0 & 0 & 0 & 0 \\ 2.2 & 8.3 & 0 & 3.7 & 1.3 & 3.8 \\ 0 & 0 & 4.2 & 0 & 0 & 0 \\ 5.4 & 0 & 0 & 9.2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1.1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 8.1 \end{pmatrix}$$

Memory footprint of COO format: nz(val) + 2*nz(int)

Memory footprint of CSR format: nz(val) + nz(int) + (n+1) (int)

CSR format

Input A, x, y Output $y = A \cdot x$

- Matrix A contains only few nonzero elements.
- Storing all entries results in large overhead (memory & computation).
- Idea: Store only nonzero elements [nz] explicitly.

Need to also store location of nonzero elements!

$$A = \begin{pmatrix} 5.4 & 1.1 & 0 & 0 & 0 & 0 \\ 2.2 & 8.3 & 0 & 3.7 & 1.3 & 3.8 \\ 0 & 0 & 4.2 & 0 & 0 & 0 \\ 5.4 & 0 & 0 & 9.2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1.1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 8.1 \end{pmatrix}$$

Memory footprint of COO format: nz(val) + 2*nz(int)

Memory footprint of CSR format: nz(val) + nz(int) + (n+1) (int)

CSR format

 $\operatorname{Input} A, x, y \quad \operatorname{Output} y = A \cdot x$

Matrix A continuous formula AStoring all ent Hands-on Exercise: Convert this matrix into CSR format: Idea: Store or Need to $\left(\begin{array}{cccccccc}
0 & 0 & 0 & 4 & 2 & 0 \\
0 & 2 & 3 & 0 & 0 & 1 \\
0 & 0 & 0 & 0 & 0 & 0 \\
1 & 2 & 3 & 4 & 3 & 1 \\
0 & 2 & 0 & 0 & 0 & 1 \\
1 & 2 & 3 & 0 & 0 & 1
\end{array}\right)$ format: ormat: $value = \begin{bmatrix} 5 \end{bmatrix}$ $colidx = \begin{bmatrix} 0 \end{bmatrix}$ Compute the memory requirement (# vals + # int) $rowptr = \begin{bmatrix} 0 & 2 & 7 & 8 \end{bmatrix}$ 10 Points to the first element in each row

Matrix A continuous formula A Storing all ent Hands-on Exercise: Convert this matrix into CSR format: Idea: Store or Need to $\left(\begin{array}{cccccccc}
0 & 0 & 0 & 4 & 2 & 0 \\
0 & 2 & 3 & 0 & 0 & 1 \\
0 & 0 & 0 & 0 & 0 & 0 \\
1 & 2 & 3 & 4 & 3 & 1 \\
0 & 2 & 0 & 0 & 0 & 1 \\
1 & 2 & 3 & 0 & 0 & 1
\end{array}\right)$ format: ormat: $value = \begin{bmatrix} 5 \end{bmatrix}$ $colidx = \begin{bmatrix} 0 \end{bmatrix}$ Compute the memory requirement (# vals + # int) 17 vals + 24 int $rowptr = \begin{bmatrix} 0 & 2 & 7 & 8 \end{bmatrix}$ 10 11 (12) Points to the first element in each row

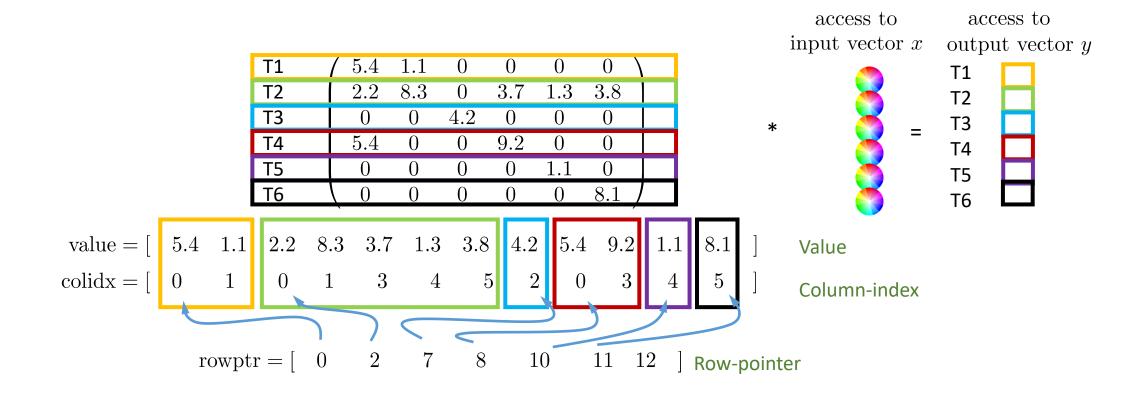
How to parallelize this?

$$A = \begin{pmatrix} 5.4 & 1.1 & 0 & 0 & 0 & 0 \\ 2.2 & 8.3 & 0 & 3.7 & 1.3 & 3.8 \\ 0 & 0 & 4.2 & 0 & 0 & 0 \\ 5.4 & 0 & 0 & 9.2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1.1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 8.1 \end{pmatrix}$$

CSR SpMV

 $\operatorname{Input} A, x, y \quad \operatorname{Output} y = A \cdot x$

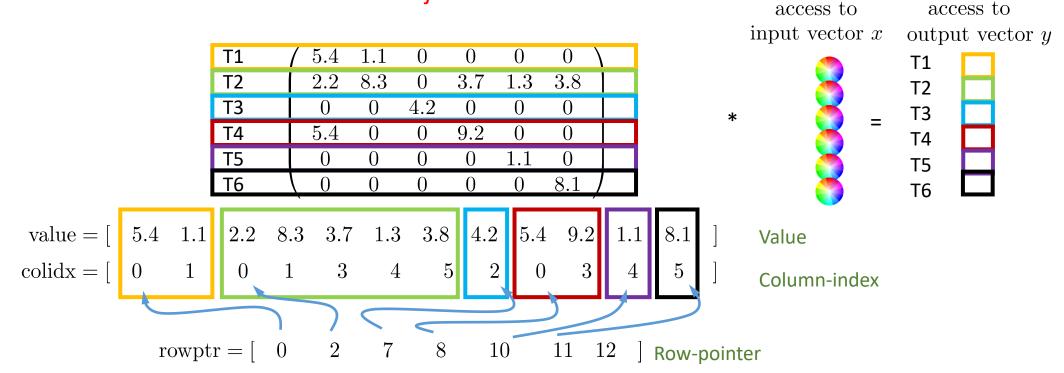
- Parallelize by rows:
 - Every "thread" handles the computation of one sum in local memory.



CSR SpMV

 $\operatorname{Input} A, x, y \quad \operatorname{Output} y = A \cdot x$

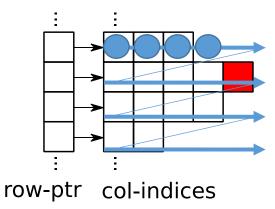
- Parallelize by rows:
 - Every "thread" handles the computation of one sum in local memory.
 - Significant workload imbalance!
 - Can not store the matrix in Col-Major format for coalesced access!



```
for( row=0; row<n; row++ )
{
    sum = 0.0;
    for( j=rowptr[row]; j<rowptr[row+1]; j++)
        sum += values[ j ] * x[ colind[j] ];
    y[ row ] = alpha * sum;
}</pre>
```

Storing values and columns in row-major.

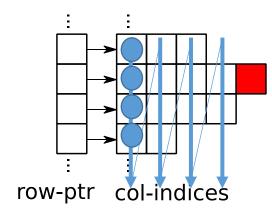
-> On GPUs: non-coalesced memory access



```
for( row=0; row<n; row++ )
{
    sum = 0.0;
    for( j=rowptr[row]; j<rowptr[row+1]; j++)
        sum += values[ j ] * x[ colind[j] ];
    y[ row ] = alpha * sum;
}</pre>
```

Storing values and columns in row-major.

-> On GPUs: non-coalesced memory access



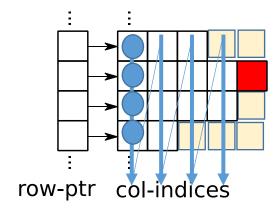
Can we use column-major?

-> Only if all rows contain the same number of nonzero elements

```
for( row=0; row<n; row++ )
{
    sum = 0.0;
    for( j=rowptr[row]; j<rowptr[row+1]; j++)
        sum += values[ j ] * x[ colind[j] ];
    y[ row ] = alpha * sum;
}</pre>
```

Storing values and columns in row-major.

-> On GPUs: non-coalesced memory access



Can we use column-major?

-> Only if all rows contain the same number of nonzero elements

 $\operatorname{Input} A, x, y \quad \operatorname{Output} y = A \cdot x$

$$A = \begin{pmatrix} 5.4 & 1.1 & 0 & 0 & 0 & 0 \\ 2.2 & 8.3 & 0 & 3.7 & 1.3 & 3.8 \\ 0 & 0 & 4.2 & 0 & 0 & 0 \\ 5.4 & 0 & 0 & 9.2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1.1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 8.1 \end{pmatrix}$$

Input A, x, y Output $y = A \cdot x$

`Left-align nonzero elements'

$$A = \begin{pmatrix} 5.4 < 1.1 & 0 & 0 & 0 & 0 \\ 2.2 < 8.3 & 0 & 3.7 & 1.3 & 3.8 \\ 0 < 0 & 4.2 & 0 & 0 & 0 \\ 5.4 < 0 & 0 & 9.2 & 0 & 0 \\ 0 < 0 & 0 & 0 & 1.1 & 0 \\ 0 < 0 & 0 & 0 & 0 & 8.1 \end{pmatrix}$$

$$\begin{bmatrix} 5.4 & 1.1 & 0 & 0 & 0 & 0 \\ 2.2 & 8.3 & 3.7 & 1.3 & 3.8 & 0 \\ 4.2 & 0 & 0 & 0 & 0 & 0 \\ 5.4 & 9.2 & 0 & 0 & 0 & 0 \\ 1.1 & 0 & 0 & 0 & 0 & 0 \\ 8.1 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0 & 1 & - & - & - & - \\ 0 & 1 & 3 & 4 & 5 & - \\ 2 & - & - & - & - & - \\ 0 & 3 & - & - & - & - \\ 4 & - & - & - & - & - \\ 5 & - & - & - & - & - \end{bmatrix}$$

$$\begin{bmatrix} 0 & 1 & - & - & - & - \\ 0 & 1 & 3 & 4 & 5 & - \\ 2 & - & - & - & - & - \\ 0 & 3 & - & - & - & - \\ 4 & - & - & - & - & - \\ 5 & - & - & - & - & - \end{bmatrix}$$

Input A, x, y Output $y = A \cdot x$

`Left-align nonzero elements'

$$A = \begin{pmatrix} 5.4 & 1.1 & 0 & 0 & 0 & 0 \\ 2.2 & 8.3 & 0 & 3.7 & 1.3 & 3.8 \\ 0 & 0 & 4.2 & 0 & 0 & 0 \\ 5.4 & 0 & 0 & 9.2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1.1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 8.1 \end{pmatrix}$$

$$\begin{bmatrix} 5.4 & 1.1 & 0 & 0 & 0 & 0 \\ 2.2 & 8.3 & 3.7 & 1.3 & 3.8 & 0 \\ 4.2 & 0 & 0 & 0 & 0 & 0 \\ 5.4 & 9.2 & 0 & 0 & 0 & 0 \\ 1.1 & 0 & 0 & 0 & 0 & 0 \\ 8.1 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 1 & - & - & - & - \\ 0 & 1 & 3 & 4 & 5 & - \\ 2 & - & - & - & - & - \\ 0 & 3 & - & - & - & - \\ 4 & - & - & - & - & - \\ 5 & - & - & - & - & - \\ \end{bmatrix}$$

Pad rows to uniform length

Memory volume:

values: max_nnz_row * num_rows

col-index: max_nnz_row * num_rows

no row pointer

Input A, x, y Output $y = A \cdot x$

$$A = \begin{pmatrix} 5.4 & 1.1 & 0 & 0 & 0 & 0 \\ 2.2 & 8.3 & 0 & 3.7 & 1.3 & 3.8 \\ 0 & 0 & 4.2 & 0 & 0 & 0 \\ 5.4 & 0 & 0 & 9.2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1.1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 8.1 \end{pmatrix}$$

$$\begin{bmatrix} 5.4 & 1.1 & 0 & 0 & 0 \\ 2.2 & 8.3 & 3.7 & 13 & 3.8 & 0 \\ 4.2 & 0 & 0 & 0 & 0 & 0 \\ 5.4 & 9.2 & 0 & 0 & 0 & 0 \\ 1.1 & 0 & 0 & 0 & 0 & 0 \\ 8.1 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 1 & - & - & - & - \\ 0 & 1 & 3 & 4 & 5 & - \\ 2 & - & - & - & - & - \\ 0 & 3 & - & - & - & - \\ 4 & - & - & - & - & - \\ 5 & - & - & - & - & - \end{bmatrix}$$

Pad rows to uniform length

Memory volume:

values: max_nnz_row * num_rows

col-index: max_nnz_row * num_rows

no row pointer

 $\operatorname{Input} A, x, y \quad \operatorname{Output} y = A \cdot x$

$$A = \begin{pmatrix} 5.4 & 1.1 & 0 & 0 & 0 & 0 \\ 2.2 & 8.3 & 0 & 3.7 & 1.3 & 3.8 \\ 0 & 0 & 4.2 & 0 & 0 & 0 \\ 5.4 & 0 & 0 & 9.2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1.1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 8.1 \end{pmatrix}$$

T1	5.4 1.1 0) 0 0	
T2	2.2 /8.3 /3.7 /	3 3.8 0	0 1 3 4 5 -
T3	4.2 / 0 / 0 /) 0 0	
T4	$oxed{5.4/9.2}$) 0 0	0 3
T5	1.1 0/ 0/) 0 0	4 -
T6	8,100) 0 0	5

Pad rows to uniform length

Memory volume:

values: max_nnz_row * num_rows

col-index: max_nnz_row * num_rows

no row pointer

Input A, x, y Output $y = A \cdot x$

$$A = \begin{pmatrix} 5.4 & 1.1 & 0 & 0 & 0 & 0 \\ 2.2 & 8.3 & 0 & 3.7 & 1.3 & 3.8 \\ 0 & 0 & 4.2 & 0 & 0 & 0 \\ 5.4 & 0 & 0 & 9.2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1.1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 8.1 \end{pmatrix}$$

T1	5.4 1.1 0	0 0	0	1 –		-
T2	2.2 /8.3 /3.7	$1 \ 3 \ 3.8 \ 0$	0	1 3	4 5	_
Т3	$oxed{4.2 / 0 / 0}$	0 0	2			_
T4	$oxed{5.4/9.2}$	0 0	0	3 –		_
T5	1.1 0/ 0/	0 0	4			_
Т6	8,1 0 0	0 0	5			_

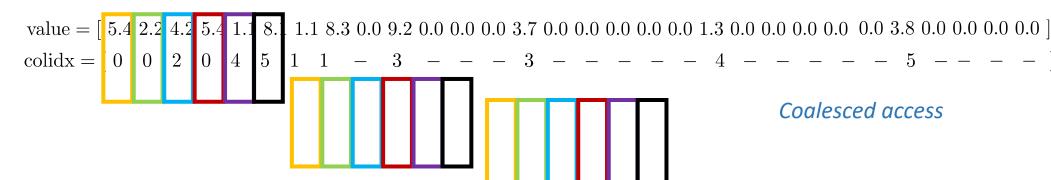
Pad rows to uniform length

Memory volume:

values: max_nnz_row * num_rows

col-index: max_nnz_row * num_rows

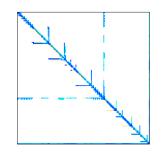
no row pointer



Performance Engineering for Sparse Linear Solvers

Input A, x, y Output $y = A \cdot x$

$$A = \begin{pmatrix} 5.4 & 1.1 & 0 & 0 & 0 & 0 \\ 2.2 & 8.3 & 0 & 3.7 & 1.3 & 3.8 \\ 0 & 0 & 4.2 & 0 & 0 & 0 \\ 5.4 & 0 & 0 & 9.2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1.1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 8.1 \end{pmatrix}$$



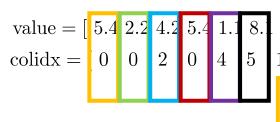
T1	5.4 1.1 0	0 0	0 1
T2	2.2 /8.3 /3.7 /1 3	3.8 0	0 1 3 4 5 -
Т3	4.2 0 / 0 / 0	0 0	
T4	5.4/9.2/0/0	0 0	0 3
T5	1.1 0/ 0/ 0/	0 0	4
T6	8,1 0 0	0 0	5

Pad rows to uniform length

Memory volume:

values: max_nnz_row * num_rows
col-index: max_nnz_row * num_rows

no row pointer



 $value = \begin{bmatrix} 5.4 & 2.2 & 4.2 & 5.4 & 1.1 & 8. & 1.1 & 8.3 & 0.0 & 9.2 & 0.0 & 0$

1 1 - 3 - - - 3 - - - - 4 - - - - 5 - - -]



Coalesced access
Can be wasteful (overhead)

Input A, x, y Output $y = A \cdot x$

$$A = \begin{pmatrix} 5.4 & 1.1 & 0 & 0 & 0 & 0 \\ 2.2 & 8.3 & 0 & 3.7 & 1.3 & 3.8 \\ 0 & 0 & 4.2 & 0 & 0 & 0 \end{pmatrix}$$

Hands-on Exercise: Convert this matrix into ELL format:

T1	5.4
T2	2.2
Т3	4.2
T4	5.4
T5	1.1
T6	8,/1

$$value = \begin{bmatrix} 5.4 & 2.2 \\ colidx = \begin{bmatrix} 0 & 0 \end{bmatrix}$$

0	0	0	4	2	0	\
0	2	3	0	0	1	
0	0	0	0	0	0	
1	2	3	4	3	1	
0	2	0	0	0	1	
1	2	3	0	0	1	

Compute the memory requirement (# vals + # int)



length

ow * num_rows
'_row * num_rows

0.0 0.0]

erhead)

Input A, x, y Output $y = A \cdot x$

$$A = \begin{pmatrix} 5.4 & 1.1 & 0 & 0 & 0 & 0 \\ 2.2 & 8.3 & 0 & 3.7 & 1.3 & 3.8 \\ 0 & 0 & 4.2 & 0 & 0 & 0 \end{pmatrix}$$

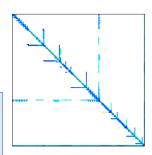
Hands-on Exercise: Convert this matrix into ELL format:

T1	5.4
T2	2.2
T3	4.2
T4	5.4
T5	1.1
T6	8,/1

$$value = \begin{bmatrix} 5.4 & 2.2 \\ colidx = \begin{bmatrix} 0 & 0 \end{bmatrix}$$

0	0	0	4	2	0	1
0	2	3		0	1	
0	0	0	0	0	0	
1	2	3	4	3	1	
0	2	0	0	0	1	
1	2	3	0	0	1	

Compute the memory requirement (# vals + # int)



length

ow * num_rows
'_row * num_rows

0.0 0.0]

erhead)

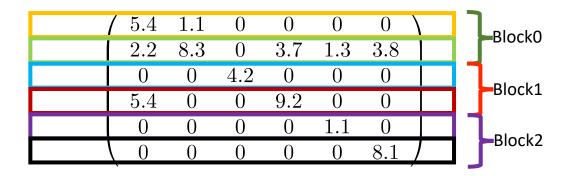
Input A, x, y Output $y = A \cdot x$

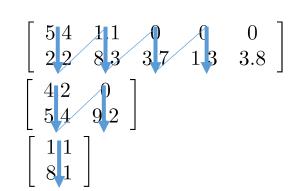
- Partition the matrix into blocks & use ELL for the distinct blocks.
 - Reduce overhead of ELL.
 - Can still store col-major.

$$A = \begin{pmatrix} 5.4 & 1.1 & 0 & 0 & 0 & 0 \\ 2.2 & 8.3 & 0 & 3.7 & 1.3 & 3.8 \\ 0 & 0 & 4.2 & 0 & 0 & 0 \\ 5.4 & 0 & 0 & 9.2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1.1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 8.1 \end{pmatrix}$$

Input A, x, y Output $y = A \cdot x$

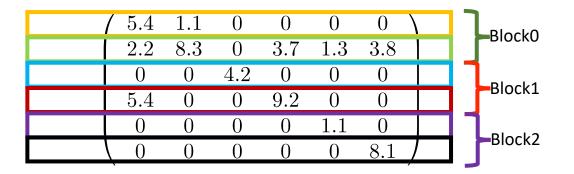
- Partition the matrix into blocks & use ELL for the distinct blocks.
 - Reduce overhead of ELL.
 - Can still store col-major.





Input A, x, y Output $y = A \cdot x$

- Partition the matrix into blocks & use ELL for the distinct blocks.
 - Reduce overhead of ELL.
 - Can still store col-major.



$$\begin{bmatrix} 5 & 4 & 1 & 1 & 0 & 0 & 0 \\ 2 & 2 & 8 & 3 & 3 & 7 & 1 & 3 & 3 & 8 \end{bmatrix}$$

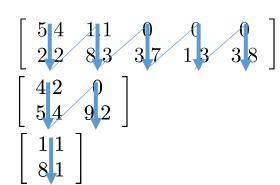
$$\begin{bmatrix} 4 & 2 & 0 \\ 5 & 4 & 9 & 2 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 1 \\ 8 & 1 \end{bmatrix}$$

Input A, x, y Output $y = A \cdot x$

- Partition the matrix into blocks & use ELL for the distinct blocks.
 - Reduce overhead of ELL.
 - Can still store col-major.
 - Need for a row pointer.

/	5.4	1.1	0	0	0	0	1		Block0
	2.2	8.3	0	3.7	1.3	3.8			Біоско
	0	0	4.2	0	0	0		_	Dlook1
	5.4	0	0	9.2	0	0			–Block1
	0	0	0	0	1.1	0			Dlask2
	0	0	0	0	0	8.1	J		Block2



sliced-ELL format:

$$value = \begin{bmatrix} 5.4 & 2.2 & 1.1 & 8.3 & 0.0 & 3.7 & 0.0 & 1.3 & 0.0 & 3.8 \\ 0 & 0 & 1 & 1 & - & 3 & - & 4 & - & 5 \end{bmatrix} \begin{bmatrix} 4.2 & 5.4 & 0.0 & 9.2 \\ 2 & 0 & - & 3 \end{bmatrix} \begin{bmatrix} 1.1 & 8.1 \\ 4 & 5 \end{bmatrix}$$

$$colidx = \begin{bmatrix} 0 & 1 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 & - & 3 & - & 4 \\ 0 & 0 & 1 & 1 & - & 3 & - & 4 \end{bmatrix}$$

$$Slice \ matrix \ into \ blocks, \ store \ blocks \ in \ ELL \ format \ with \ offset-pointer.$$

Hands-on Exercise: Convert this matrix into Sliced-ELL format (SELL-2):

$$\left(\begin{array}{ccccccc}
0 & 0 & 0 & 4 & 2 & 0 \\
0 & 2 & 3 & 0 & 0 & 1 \\
0 & 0 & 0 & 0 & 0 & 0 \\
1 & 2 & 3 & 4 & 3 & 1 \\
0 & 2 & 0 & 0 & 0 & 1 \\
1 & 2 & 3 & 0 & 0 & 1
\end{array}\right)$$

Compute the memory requirement (# vals + # int)

Hands-on Exercise: Convert this matrix into Sliced-ELL format (SELL-2):

$$\left(\begin{array}{ccccccc}
0 & 0 & 0 & 4 & 2 & 0 \\
0 & 2 & 3 & 0 & 0 & 1 \\
0 & 0 & 0 & 0 & 0 & 0 \\
1 & 2 & 3 & 4 & 3 & 1 \\
0 & 2 & 0 & 0 & 0 & 1 \\
1 & 2 & 3 & 0 & 0 & 1
\end{array}\right)$$

Rowptr: 0 6 18 26

Compute the memory requirement (# vals + # int)

26 vals + 26 + 4 int

- How can we optimize this? Minimize the overhead? What is the overhead dependent on?
 - Bring rows with similar number of nonzero elements into the same block.
 - Sort rows by "length" and reorder the matrix, then convert to Sliced-ELL

Software and High-Performance Computing

A Unified Sparse Matrix Data Format for Efficient General Sparse Matrix-Vector Multiplication on Modern Processors with Wide SIMD Units

Authors: Moritz Kreutzer, Georg Hager, Gerhard Wellein, Holger Fehske, and Alan R. Bishop

<u>AUTHORS INFO & AFFILIATIONS</u>

https://doi.org/10.1137/130930352

- What happens for block-size 1?
- What happens for block size n (matrix size)?

SpMV Formats and Kernels

 $\operatorname{Input} A, x, y \quad \operatorname{Output} y = A \cdot x$

"Different kernels optimal for different problems"

COO

- can compensate workload imbalance for irregular patterns
- Efficient for MIMD processing
- Strong support for atomics needed

CSR

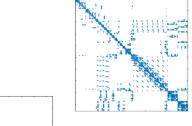
- small memory footprint
- Needs some logic for row-parallel processing
- Efficient for MIMD processing

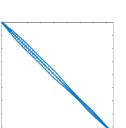
ELL

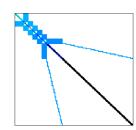
- Efficient for balanced matrices
- Enables col-major storage
- Efficient for SIMD processing

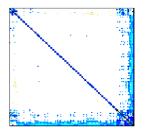
SELL-c

- Enables col-major storage
- Tunable between CSR and ELL





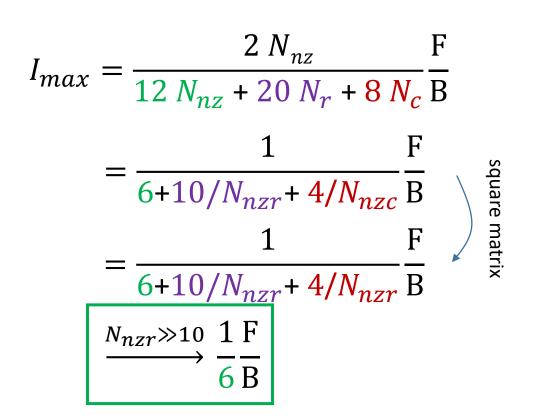


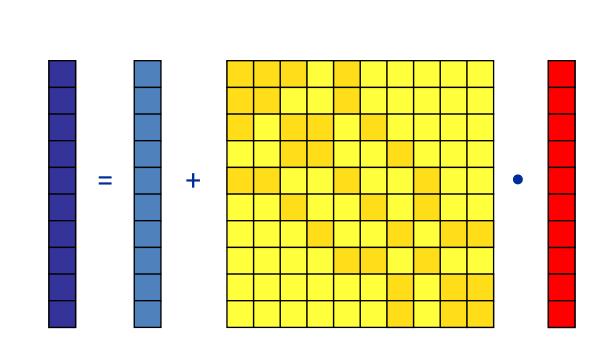


SpMV (CSR) Roofline

Optimistic intensity:

```
for (int row = 0; row < num_rows; ++row) {
  double sum = 0.0;
  for (int k = row_ptrs[row]; k < row_ptrs[row + 1]; ++k)
    sum += mat_values[k] * b[col_idxs[k]];
  x[row] += sum;
}</pre>
```





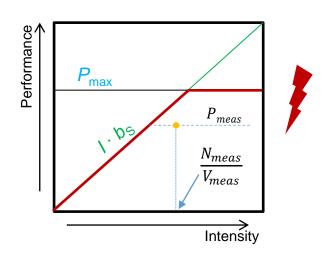
Roofline "failure" with SpMV

Reasons for performance not attaining the limit

- 1. Intensity lower than the minimum
 - More RHS traffic than the optimistic limit ($\frac{4}{N_{nzr}}$ B/F)

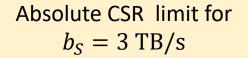


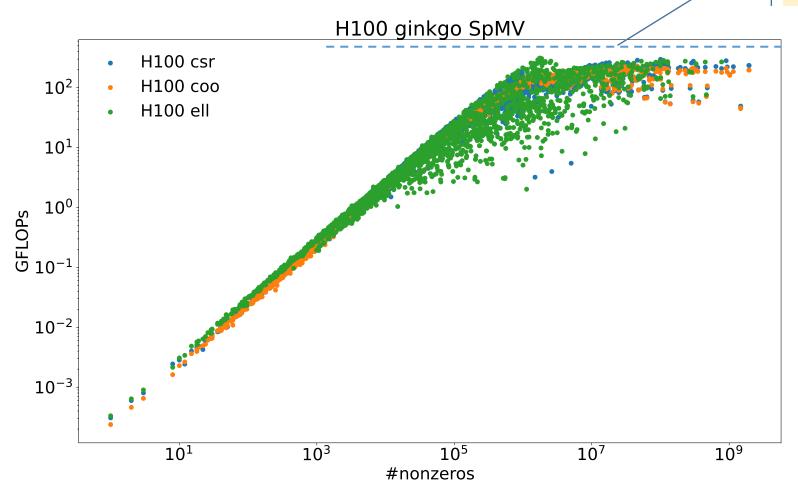
- "invisible" performance ceiling due to inefficient instructions or inefficient execution
- 3. Load imbalance
 - A single process/thread cannot saturate the memory bandwidth
- 4. Erratic memory access patterns for RHS
 - Latency dominates



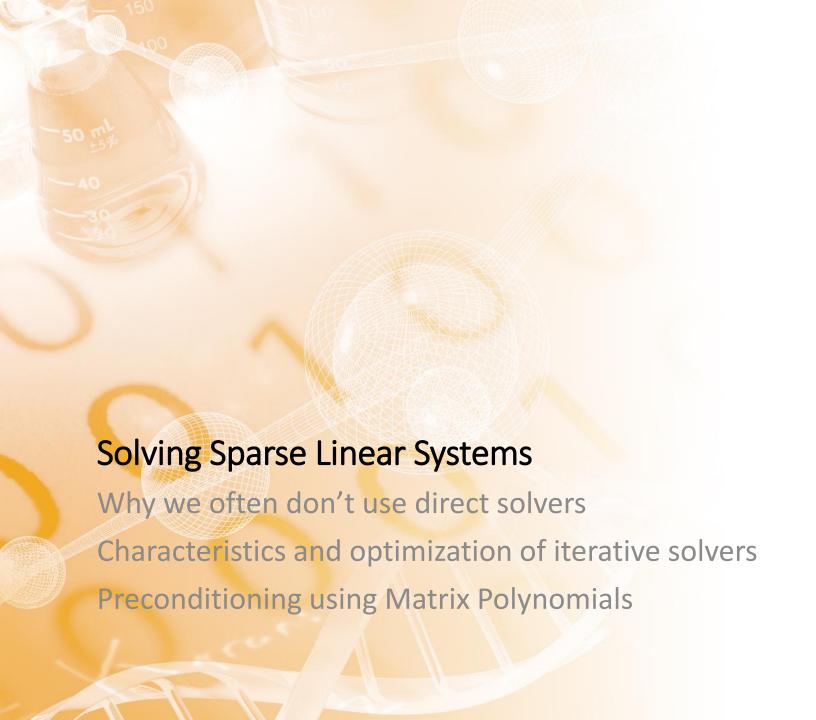
Experiences with SpMV on GPUs

Looking at ~3,000 test matrices from Suite Sparse Matrix Collection



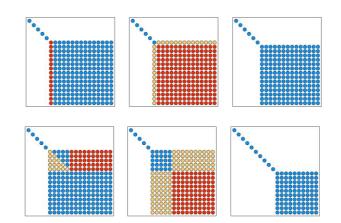






Can we use direct solvers for solving sparse problems?

$$A = \begin{pmatrix} 5.4 & 1.1 & 0 & 0 & 0 & 0 \\ 2.2 & 8.3 & 0 & 3.7 & 1.3 & 3.8 \\ 0 & 0 & 4.2 & 0 & 0 & 0 \\ 5.4 & 0 & 0 & 9.2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1.1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 8.1 \end{pmatrix}$$



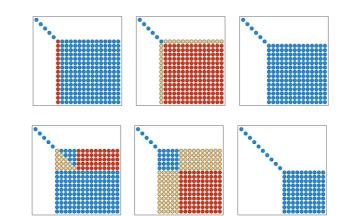
Scalar or block-LU?

Are zeros preserved in the factorization?

Can we store the fill-in?

Can we use direct solvers for solving sparse problems?

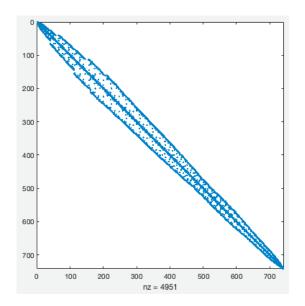
$$A = \begin{pmatrix} 5.4 & 1.1 & 0 & 0 & 0 & 0 \\ 2.2 & 8.3 & 0 & 3.7 & 1.3 & 3.8 \\ 0 & 0 & 4.2 & 0 & 0 & 0 \\ 5.4 & 0 & 0 & 9.2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1.1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 8.1 \end{pmatrix}$$



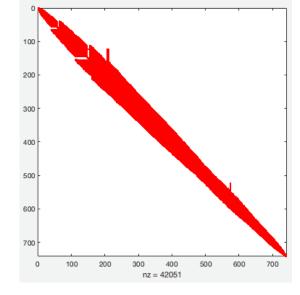
Scalar or block-LU?

Are zeros preserved in the factorization?

Can we store the fill-in?







Generate a sequence of solution approximations with increasing approximation quality.

$$x^0 \rightsquigarrow x^1 \rightsquigarrow x^2 \rightsquigarrow x^3 \rightsquigarrow \cdots$$

Generate a sequence of solution approximations with increasing approximation quality.

$$x^0 \rightsquigarrow x^1 \rightsquigarrow x^2 \rightsquigarrow x^3 \rightsquigarrow \cdots$$

Relaxations

- Base on matrix splitting
- Jacobi relaxation:

$$Ax = b$$

$$(L+D+U)x = b$$

$$Dx = b - (L+U)x$$

$$x = D^{-1}b - D^{-1}(L+U)x$$

$$x^{k+1} = D^{-1}b - D^{-1}(A-D)x^k$$

- Matrix-Vector Prod. & Vector Ops
- Low arithmetic intensity

Generate a sequence of solution approximations with increasing approximation quality.

$x^0 \rightsquigarrow x^1 \rightsquigarrow x^2 \rightsquigarrow x^3 \rightsquigarrow \cdots$

Relaxations

- Base on matrix splitting
- Jacobi relaxation:

$$Ax = b$$

$$(L+D+U)x = b$$

$$Dx = b - (L+U)x$$

$$x = D^{-1}b - D^{-1}(L+U)x$$

$$x^{k+1} = D^{-1}b - D^{-1}(A-D)x^k$$

- Matrix-Vector Prod. & Vector Ops
- Low arithmetic intensity

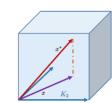
Krylov Subspace Methods

Iteratively grow Krylov subspace

$$K_i(A,r) = span\left\{r,Ar,A^2r,\dots,A^{i-1}r\right\}$$

$$K_0 \subset K_1 \subset K_2 \subset \cdots \mathbb{R}^n$$

 Approximate solution in Krylov Subspace



- Matrix-Vector Prod. & Vector Ops
- Low arithmetic intensity

Generate a sequence of solution approximations with increasing approximation quality.

Relaxations

- Base on matrix splitting
- Jacobi relaxation:

$$Ax = b$$

$$(L+D+U)x = b$$

$$Dx = b - (L+U)x$$

$$x = D^{-1}b - D^{-1}(L+U)x$$

$$x^{k+1} = D^{-1}b - D^{-1}(A-D)x^k$$

- Matrix-Vector Prod. & Vector Ops
- Low arithmetic intensity

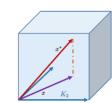
Krylov Subspace Methods

Iteratively grow Krylov subspace

$$K_i(A,r) = span\{r, Ar, A^2r, \dots, A^{i-1}r\}$$

$$K_0 \subset K_1 \subset K_2 \subset \cdots \mathbb{R}^n$$

 Approximate solution in Krylov Subspace

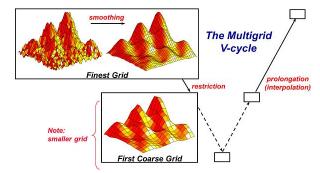


- Matrix-Vector Prod. & Vector Ops
- · Low arithmetic intensity

$$x^0 \rightsquigarrow x^1 \rightsquigarrow x^2 \rightsquigarrow x^3 \rightsquigarrow \cdots$$

Multigrid Methods

 Recursively project problem to coarser grid and solve on coarser grid



- Matrix-Vector Prod. & Vector Ops
- Low arithmetic intensity

Generate a sequence of solution approximations with increasing approximation quality.

Relaxations

- Base on matrix splitting
- Jacobi relaxation:

$$Ax = b$$

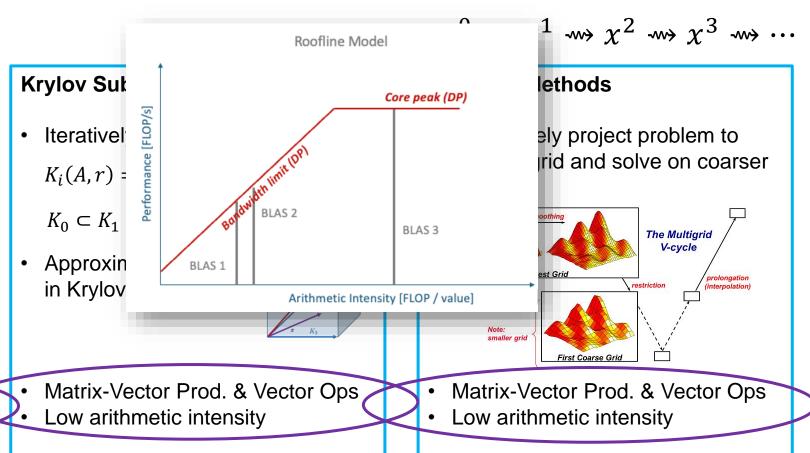
$$(L+D+U)x = b$$

$$Dx = b - (L+U)x$$

$$x = D^{-1}b - D^{-1}(L+U)x$$

$$x^{k+1} = D^{-1}b - D^{-1}(A-D)x^k$$

- Matrix-Vector Prod. & Vector Ops
- Low arithmetic intensity



Example: Conjugate Gradients (CG)

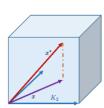
Krylov Subspace Methods

Iteratively grow Krylov subspace

$$K_i(A,r) = span\left\{r,Ar,A^2r,\dots,A^{i-1}r\right\}$$

$$K_0 \subset K_1 \subset K_2 \subset \cdots \mathbb{R}^n$$

Approximate solution in Krylov Subspace



- Matrix-Vector Prod. & Vector Ops
- Low arithmetic intensity

Optimal Krylov solver for symmetric and positive definite (SPD) matrices

Requires storing only four additional vectors

input:
$$A, b, x_0, it_{max}$$
 $r_0 = b - Ax_0$
 $p_0 = r_0$
for $k = 0, \dots, it_{max}$ do
$$\alpha_k = \frac{r_k^T r_k}{p_k^T A p_k}$$

$$x_{k+1} = x_k + \alpha_k p_k$$

$$r_{k+1} = r_k - \alpha_k A p_k$$

$$\beta_k = \frac{r_{k+1}^T r_{k+1}}{r_k^T r_k}$$

$$p_{k+1} = r_{k+1} + \beta_k p_k$$
end for

Preconditioning Iterative Solvers

Transform linear problem by multiplying both sides with $P \approx A^{-1}$ such that iterations converge faster.

$$Ax = b \Leftrightarrow PAx = Pb \Leftrightarrow \tilde{A}x = \tilde{b}$$

Transform linear problem by multiplying both sides with $P \approx A^{-1}$ such that iterations converge faster.

$Ax = b \Leftrightarrow \underbrace{PAx = Pb}_{\tilde{A}} \Leftrightarrow \tilde{A}x = \tilde{b}$

Iterative solver as preconditioner

- Multigrid
- Jacobi $D^{-1}Ax = D^{-1}b$
- Block-Jacobi



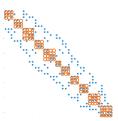
- Sparse Approximate Inverses
- Matrix-Vector Prod. & Vector Ops
- Low arithmetic intensity

Transform linear problem by multiplying both sides with $P \approx A^{-1}$ such that iterations converge faster.

$$Ax = b \Leftrightarrow PAx = Pb \Leftrightarrow \tilde{A}x = \tilde{b}$$

Iterative solver as preconditioner

- Multigrid
- Jacobi $D^{-1}Ax = D^{-1}b$
- Block-Jacobi



- Sparse Approximate Inverses
- Matrix-Vector Prod. & Vector Ops
- Low arithmetic intensity

Incomplete Factorizations

Compute LU factorization with restricted fill-in

 Replace triangular solver with iteratively solving factors

$$L \cdot y = b$$
 $U \cdot x = y$

- Matrix-Vector Prod. & Vector Ops
- Low arithmetic intensity

Transform linear problem by multiplying both sides with $P \approx A^{-1}$ such that iterations converge faster.

$$Ax = b \Leftrightarrow PAx = Pb \Leftrightarrow \tilde{A}x = \tilde{b}$$

Iterative solver as preconditioner

- Multigrid
- Jacobi $D^{-1}Ax = D^{-1}b$
- Block-Jacobi



- Sparse Approximate Inverses
- Matrix-Vector Prod. & Vector Ops
- Low arithmetic intensity

Incomplete Factorizations

Compute LU factorization with restricted fill-in

 Replace triangular solver with iteratively solving factors

$$L \cdot y = b$$
 $U \cdot x = y$

- Matrix-Vector Prod. & Vector Ops
- · Low arithmetic intensity

Polynomial preconditioners

• Choose A = M - N

$$P = \left(\sum_{i=0}^{p-1} (I - M^{-1}A)^i\right) M^{-1}$$

- Matrix-Vector Prod. & Vector Ops
- Low arithmetic intensity

Transform linear problem by multiplying both sides with $P \approx A^{-1}$ such that iterations converge faster.

$$Ax = b \Leftrightarrow \underbrace{PAx = Pb}_{\tilde{A}} \Leftrightarrow \tilde{A}x = \tilde{b}$$

preconditioners

 $\int_{0}^{1} (I - M^{-1}A)^{i} M^{-1}$

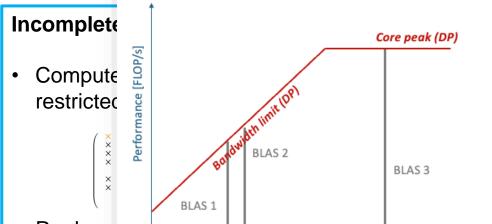
A = M - N

Iterative solver as preconditioner

- Multigrid
- Jacobi $D^{-1}Ax = D^{-1}b$
- Block-Jacobi



- Sparse Approximate Inverses
- Matrix-Vector Prod. & Vector Ops
- Low arithmetic intensity



Roofline Model

Replace
 iteratively solving factors

Arithmetic Intensity [FLOP / value]

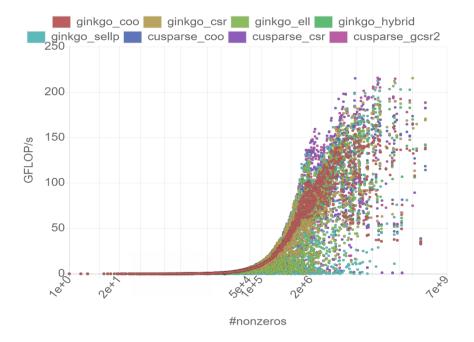
$$L \cdot y = b \quad U \cdot x = y$$

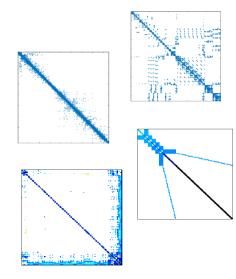
- Matrix-Vector Prod. & Vector Ops
- Low arithmetic intensity

- Matrix-Vector Prod. & Vector Ops
- Low arithmetic intensity

Optimizing Iterative Solvers & Preconditioners

- 1. Optimizing the matrix vector product as common building block
 - Optimization of sparse data format and processing scheme





Optimizing Iterative Solvers & Preconditioners

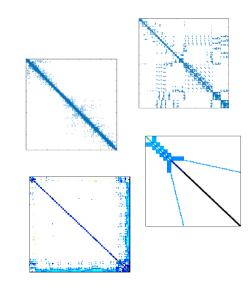
- Optimizing the matrix vector product as common building block
 - Optimization of sparse data format and processing scheme

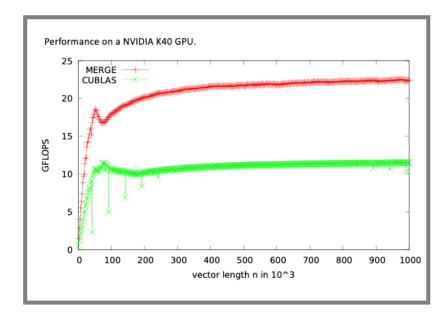
2. Cache-Aware implementation

Merging of Operations into super-kernels to reduce the memory access

```
BiCGStab Krylov solver (van der Vorst, 1992)
       1. r_0 = b - Ax_0
       2. Choose an arbitrary vector \hat{r}_0 such that (\hat{r}_0, r_0) \neq 0, e.g., \hat{r}_0 = r_0
       3. \rho_0 = \alpha = \omega_0 = 1
       4. v_0 = p_0 = 0
      5. For i = 1, 2, 3, ...
                 1. \rho_i = (\hat{r}_0, r_{i-1})
                 2. \beta = (\rho_i/\rho_{i-1})(\alpha/\omega_{i-1})
                 3. \mathbf{p}_i = \mathbf{r}_{i-1} + \beta(\mathbf{p}_{i-1} - \omega_{i-1}\mathbf{v}_{i-1})
                  4. v_i = Ap_i
                  5. \alpha = \rho_i/(\hat{\mathbf{r}}_0, \mathbf{v}_i)
                  6. \mathbf{s} = \mathbf{r}_{i-1} - \alpha \mathbf{v}_i
                  7. t = As
                  8. \omega_i = (t, s)/(t, t)
                  9. \mathbf{x}_i = \mathbf{x}_{i-1} + \alpha \mathbf{p}_i + \omega_i \mathbf{s}
                10. If x_i is accurate enough then quit
                                                                                                                     mance Engineering for Sparse Linear Solvers
                11. \mathbf{r}_i = \mathbf{s} - \omega_i \mathbf{t}
```

```
p_k := r_{k-1} + \beta (p_{k-1} - \omega_{k-1} v_{k-1})
   cuBLAS
  cublasDscal(n, beta, p, 1);
  cublasDaxpy( n, omega * beta, v, 1, p, 1);
  cublasDaxpy( n, 1.0, r, 1, p, 1 );
    3 kernels - 5n reads, 3n writes
   merge in one kernel
  p update(int n, double beta, double omega,
               double *v, double *r, double *p ){
    int i = blockldx.x * blockDim.x + threadldx.x;
    if(i<n)
      p[i] = r[i] + beta * (p[i]-omega*v[i]);
    1 kernel - 3n reads, 1n writes
```





Optimizing Iterative Solvers & Preconditioners

1. Optimizing the matrix vector product as common building block

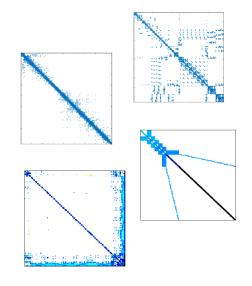
Optimization of sparse data format and processing scheme

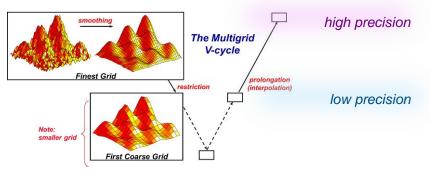
2. Cache-Aware implementation

Merging of Operations into super-kernels to reduce the memory access

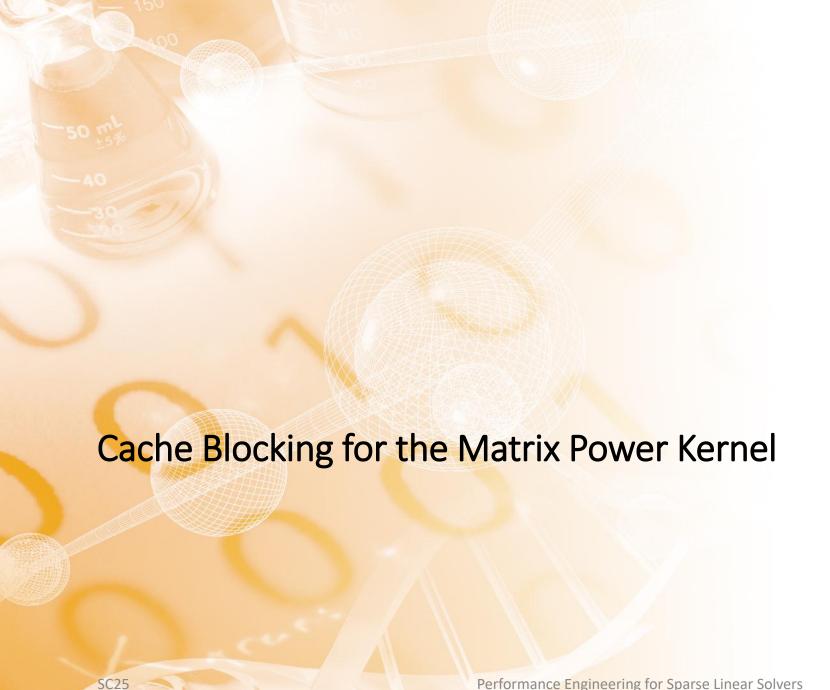
3. Replace memory access with additional computations

- Mixed Precision algorithms using low precision in parts of the computations
- Matrix Powers Kernel and cache blocking



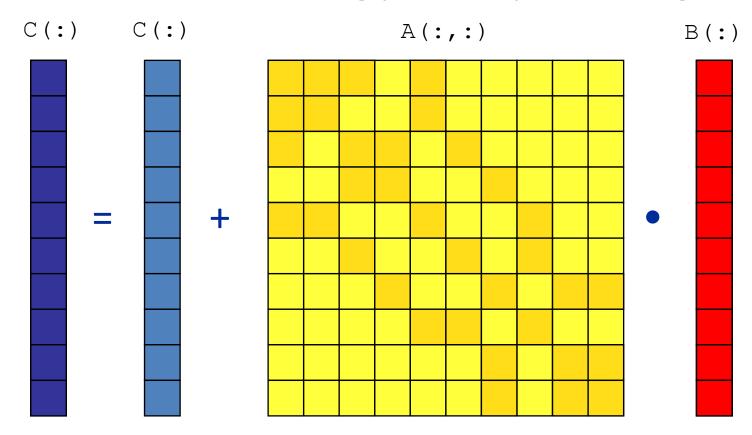






Motivation – Sparse Matrix Vector Multiplication

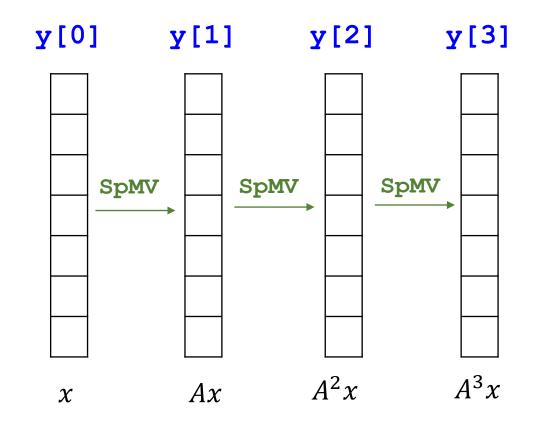
- Easy to parallelize but sparse irregular data structures / accesses
- SpMV Performance ← → Strongly Memory Bound (high code balance)



Motivation – Matrix power kernel (MPK)

- Calculate: $y = A^p x$
- Repeatedly perform back to back SpMVs

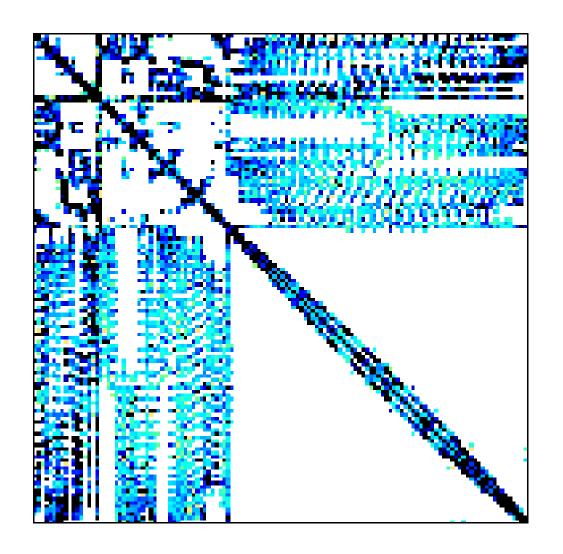
```
for k=1:p; do
  y[k] = SpMV(A, y[k-1])
done
```

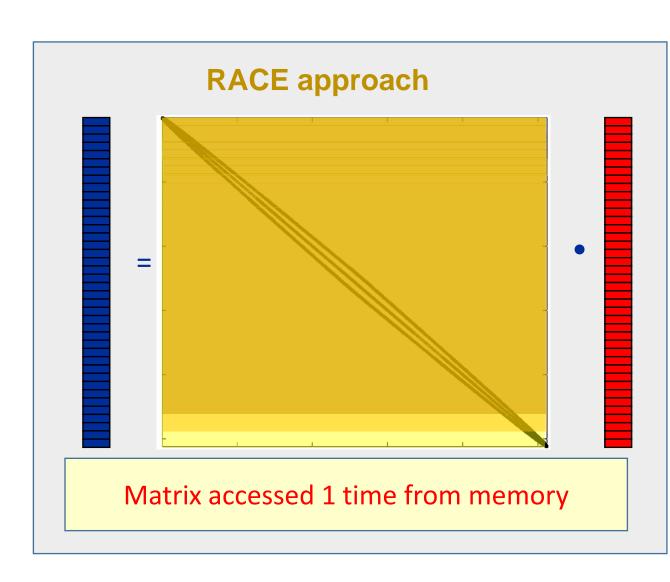


Same matrix A loaded p times from main memory!!!

How to cache the matrix A across the matrix power calculation?

Matrix power – Traditional approach vs. Cache Blocking





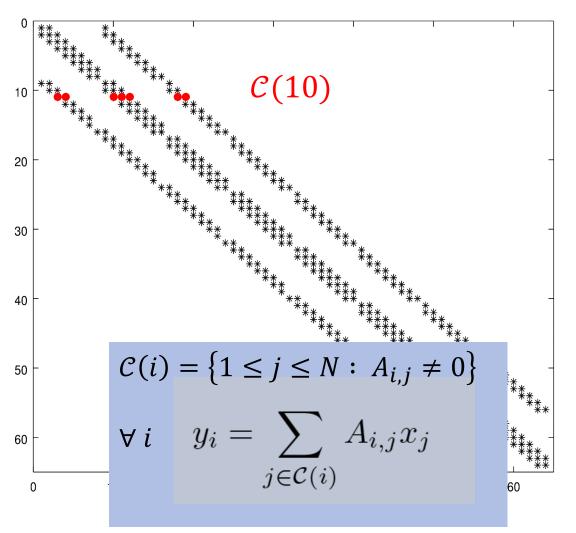
How to do that in general for sparse matrices?

SpMV – Graph Traversal – RACE

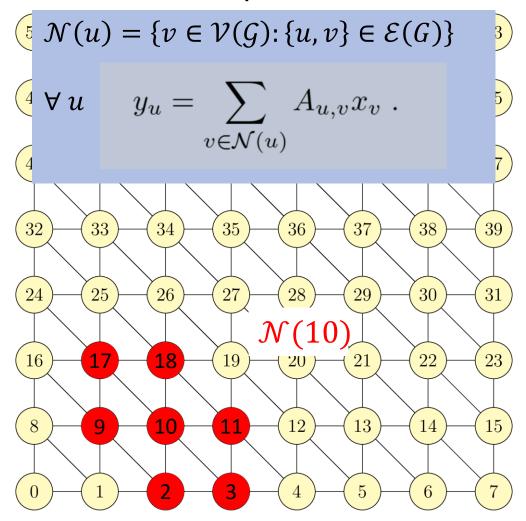


Sample matrix and its graph representation

Symmetric Matrix

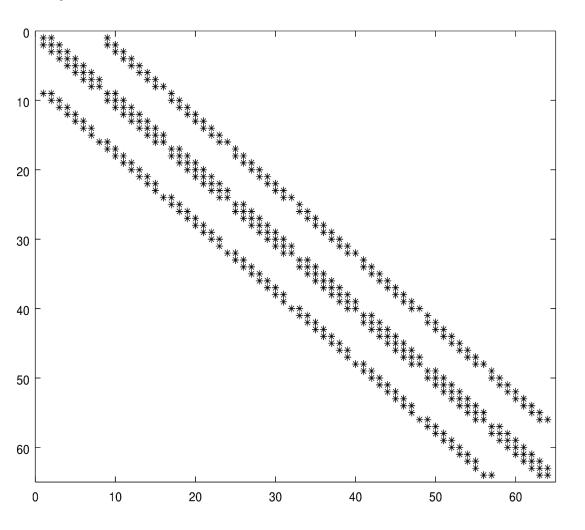


Undirected Graph

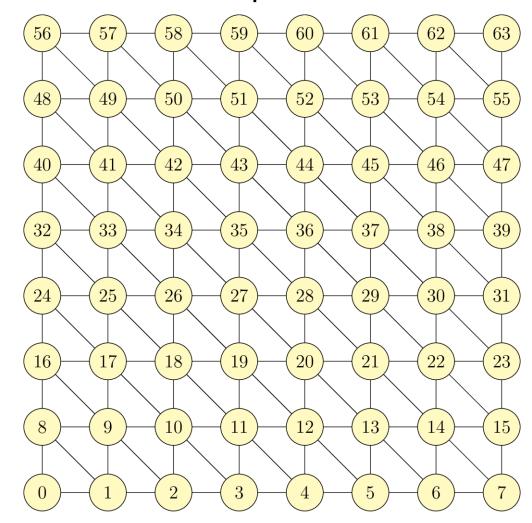


Sample matrix and its graph representation

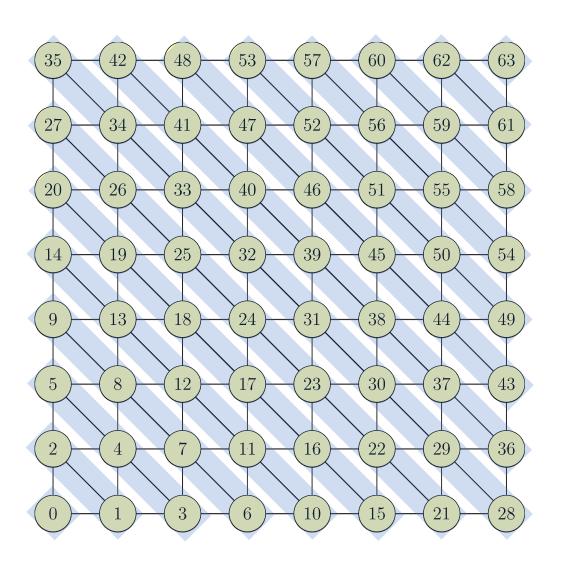
Symmetric Matrix

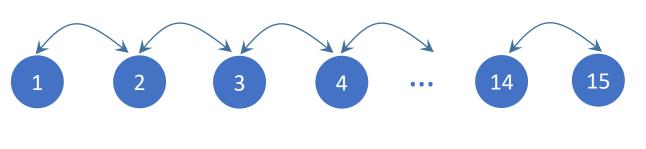


Undirected Graph



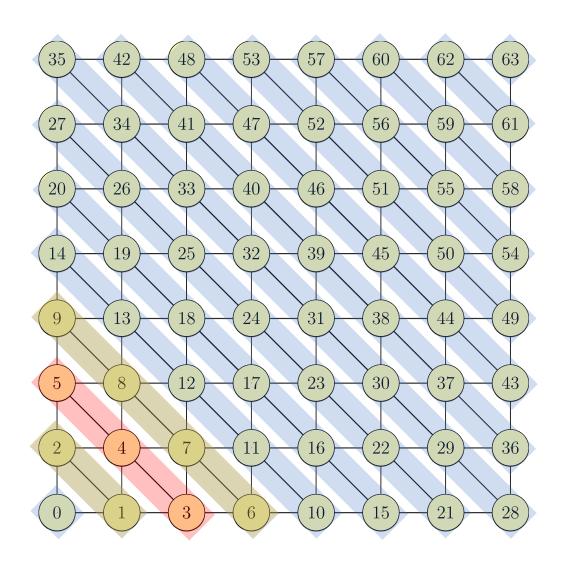
RACE

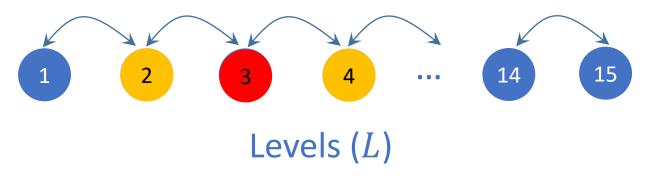




levels

RACE





Key property
$$\mathcal{N}\big(L(i)\big) = \{L(i-1), L(i), L(i+1)\}$$

 $A^p x$ computations on L(i) will require $A^{p-1} x$ to be complete on L(i-1), L(i), L(i+1)



do k = 1, p
y(:, k) = SpMV(A, y(:,k-1))
enddo

No cache blocking!

Levels

- •

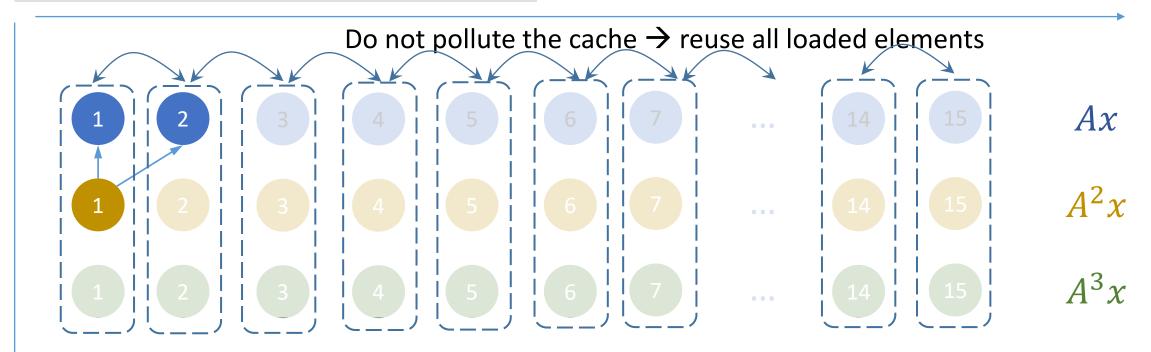
Ax

 A^2x

 A^3x

```
do k = 1, p
y(:, k) = SpMV(A, y(:,k-1))
enddo
```

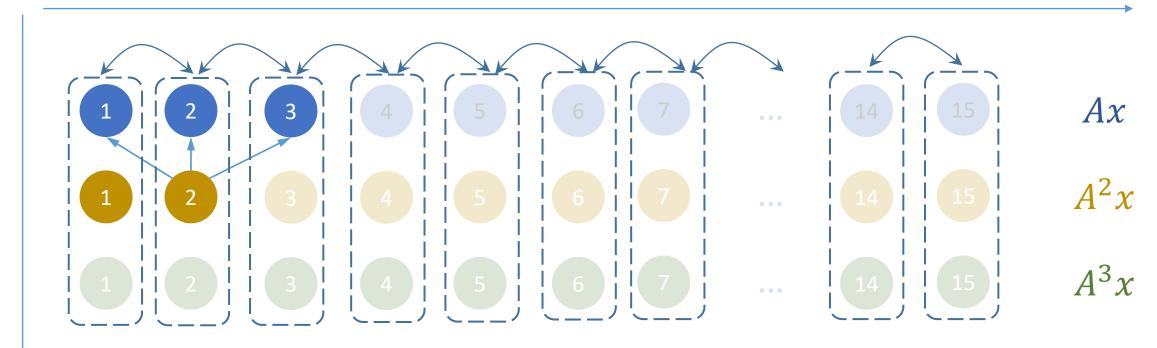
Levels



When updating level 1, indirect reads also go to level 2

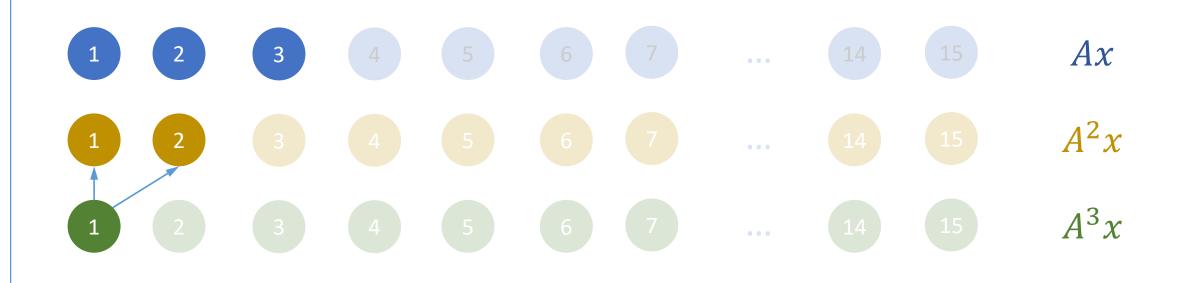
```
do k = 1, p
y(:, k) = SpMV(A, y(:,k-1))
enddo
```

Levels

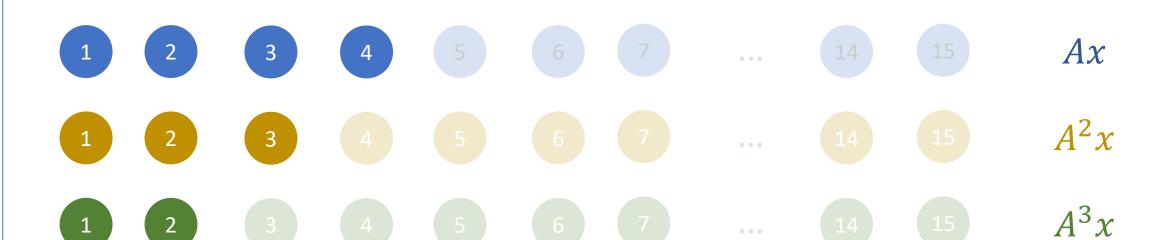


When updating level 2, indirect reads also go to levels 1 and 3

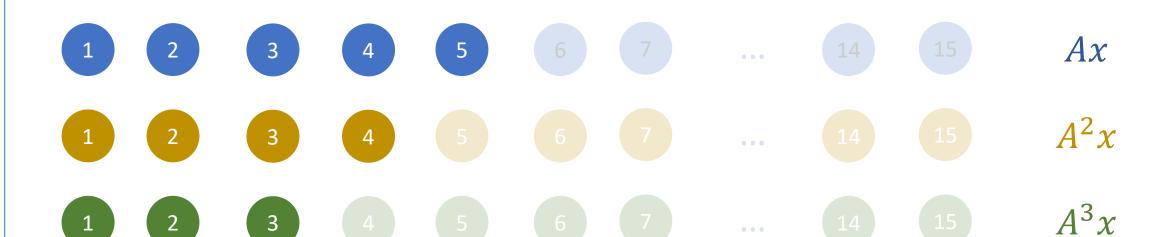
```
do k = 1, p
y(:, k) = SpMV(A, y(:,k-1))
enddo
```



```
do k = 1, p
y(:, k) = SpMV(A, y(:,k-1))
enddo
```



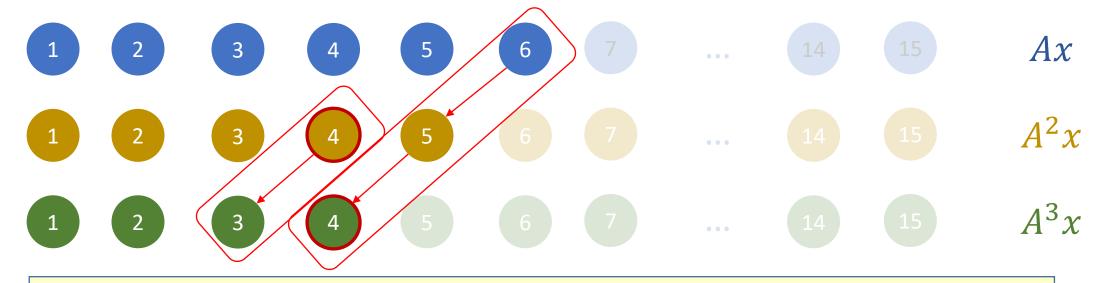
```
do k = 1, p
y(:, k) = SpMV(A, y(:,k-1))
enddo
```



RACE: MPK implementation idea



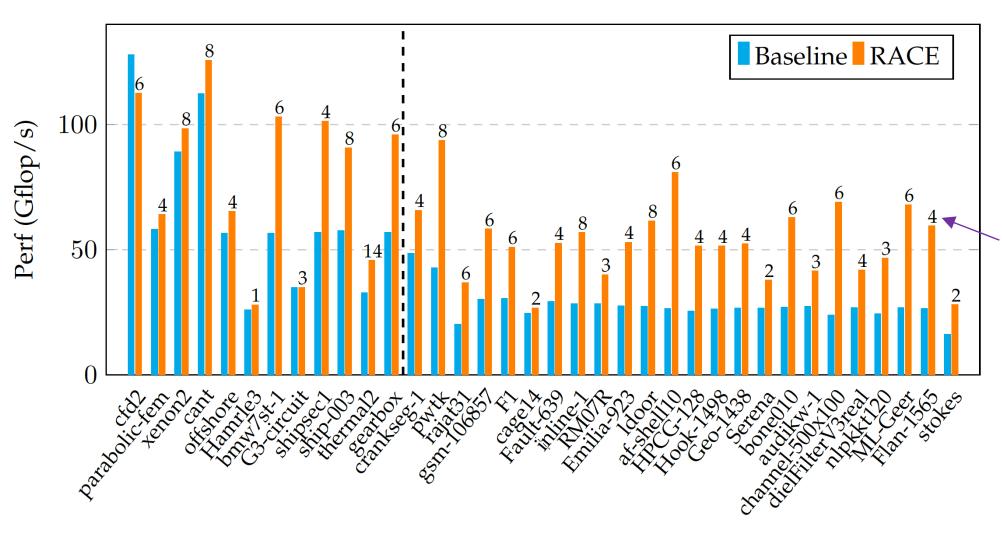
RACE – Input parameters and its influence



$$(p+1) \times N_{nz}(L) \times 12 \text{ bytes} < C$$

$$N_{nz}(L)$$
 — avg. non-zeros in a level C — cache size

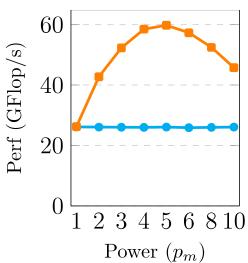
Matrix power kernel: Performance – Intel Ice Lake



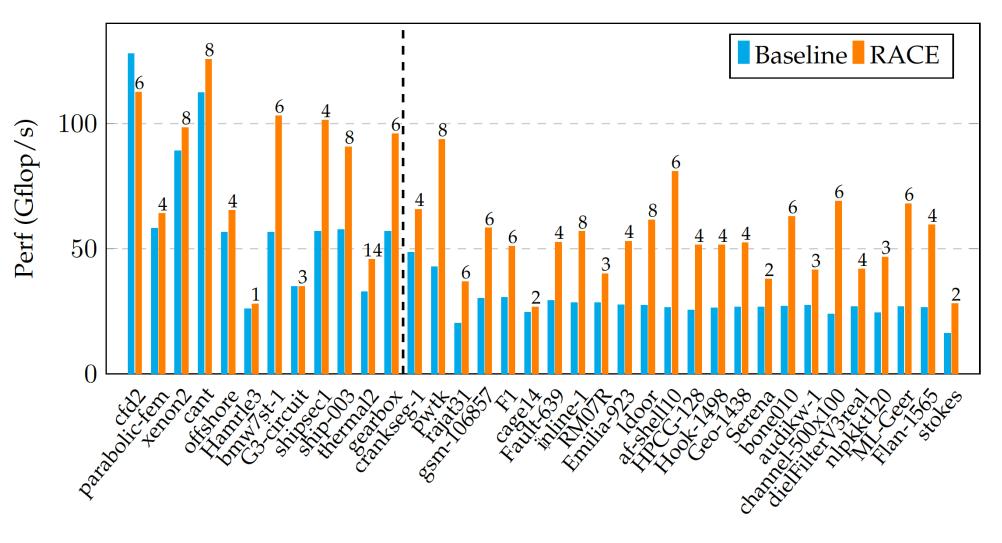
Intel Xeon Platinum 8368 (Ice Lake)

- 38 cores
- 104 MB cache (L2+L3)

Power value with maximum performance.



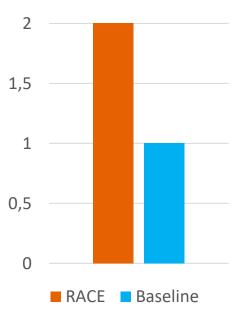
Matrix power kernel: Performance – Intel Ice Lake



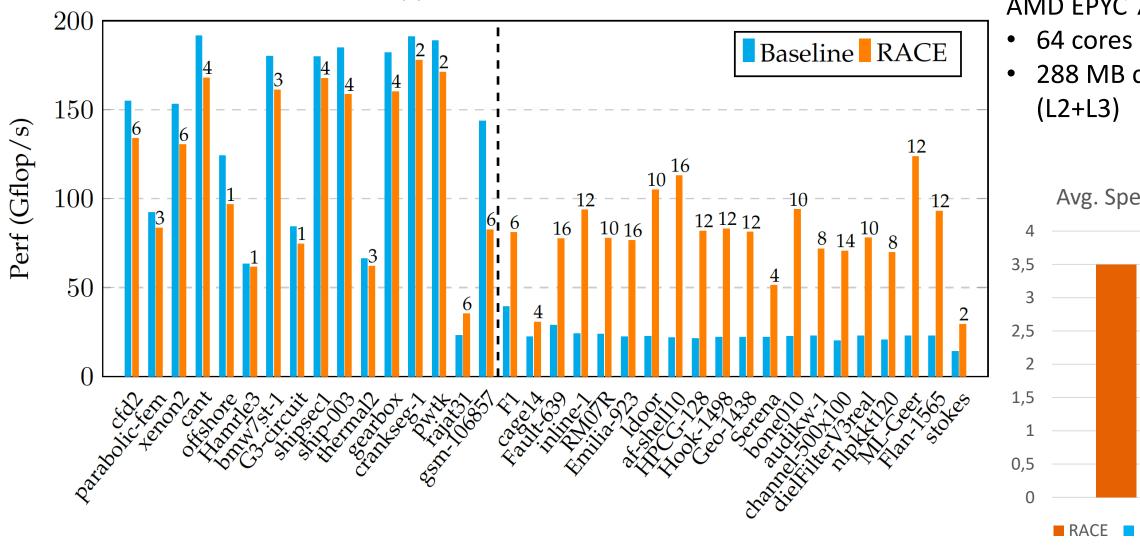
Intel Xeon Platinum 8368 (Ice Lake)

- 38 cores
- 104 MB cache (L2+L3)



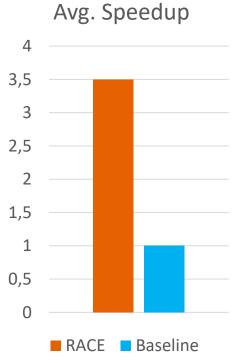


Matrix power kernel: Performance – AMD Rome



AMD EPYC 7662

288 MB cache



RACE - summary

- Inner kernel: OpenMP parallel standard SpMV routine
- Overhead: BFS & Set up of data structures (approx. ≤ 50 SpMVs)
- Parameters: Power (p_m), Available Cache Size, Max. recursion depth
- Cache size $\leftarrow \rightarrow$ max. polynomial degree (p_m)
 - Larger caches \rightarrow larger $p_m \rightarrow$ better performance
 - Polynomial degree higher than $p_m \rightarrow$ Computation in chunks of p_m
- No loss of accuracy!

RACE – MPK applications

- Exponential Integrators → Polynomial approximations
- s-step Krylov methods (CA-GMRES)
- Polynomial preconditioning
- Algebraic Multigrid smoothers

International Journal of HIGH PERFORMANCE COMPUTING APPLICATIONS

Research Paper

Algebraic temporal blocking for sparse iterative solvers on multi-core CPUs

Christie Alappat¹, Jonas Thies², Georg Hager¹, Holger Fehske¹ and Gerhard Wellein^{1,2,3}

The International Journal of High Performance Computing Applications 2024, Vol. 0(0) 1–21 © The Author(s) 2024



Article reuse guidelines: sagepub.com/journals-permissions DOI: 10.1177/10943420241283828 journals.sagepub.com/home/hpc



https://doi.org/10.1177/10943420241283828



Using RACE

```
# include <RACE/interface.h>
RACE::dist k = RACE::POWER;
Nt = omp get num threads();
RACE::Interface race (Nr, Nt, k, rowPtr, col);
//power value; here 4
int pm = 4;
//cache size in bytes; here 30 MB
double C = 30*1024*1024;
//perform pre-processing, find levels
race.RACEColor(pm, C);
int *perm, * invPerm , permLen=Nr;
race.getPerm(&perm, &permLen);
race.getInvPerm(&invPerm, &permLen);
//permute matrix and vector data structures
permute (perm, invPerm);
```

```
Pre-processing
```

```
struct functionArg
  //user-defined struct for input and output
  //arguments of the call-back function
  int Nr;
//user-defined call-back function
void foo(int row s, int row e, int pow, void * voidArg)
 functionArg * arg = (functionArg *) voidArg;
functionArg* args = new functionArg;
//fill args
args->Nr = 1000;
. . .
void* voidArgs = (void*) args;
int foo id = race.registerFunction(&foo, voidArgs, pm);
race.executeFunction(foo id);
```

Processing

Neumann polynomial apply

$$w = (I-L)^k A (I-U)^k v$$
 Cache blocking
$$t_2 = A t_1$$

$$w = (I-L)^k t_2$$
 Cache blocking

Can we do better?

Neumann polynomial apply

$$w = (I-L)^k A (I-U)^k v$$

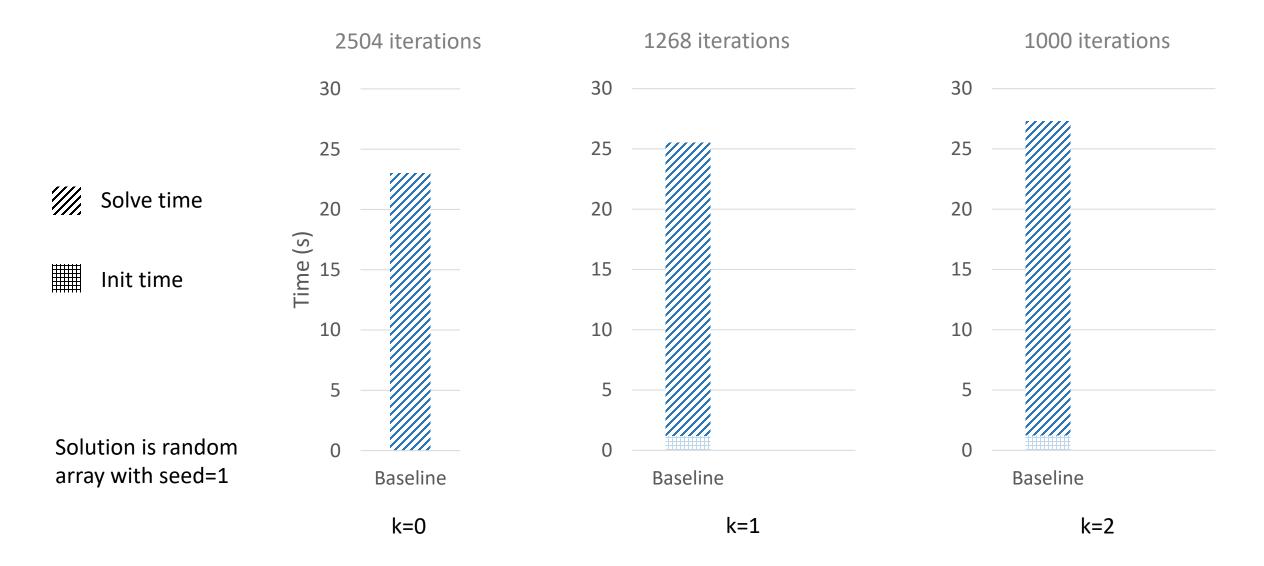
$$t_1 = (I-U)^k v$$

$$t_2 = At_1 = (L+U)t_1$$

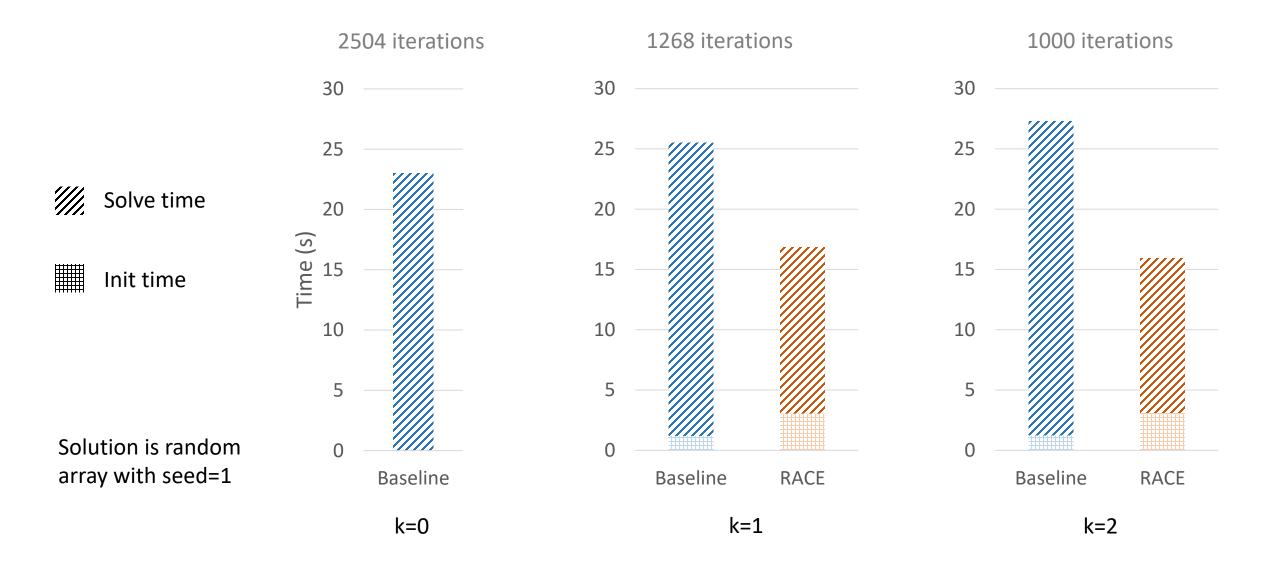
$$w = (I-L)^k t_2$$
Cache blocking

Total power = 2k + 1

Time to solve Laplace2000x2000 to 1e-3 tolerance on 1 NUMA domain (18c) of Intel Ice lake (Fritz)



Time to solve Laplace2000x2000 to 1e-3 tolerance on 1 NUMA domain (18c) of Intel Ice lake (Fritz)





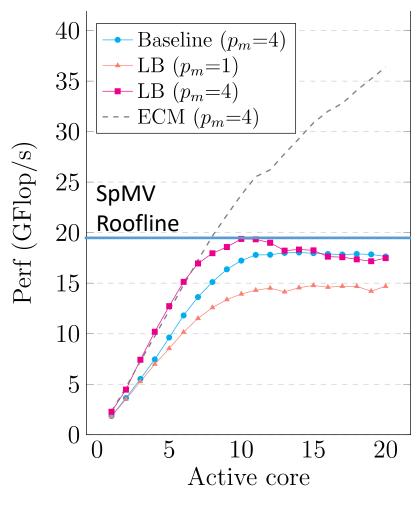
MPK – existing caching approaches

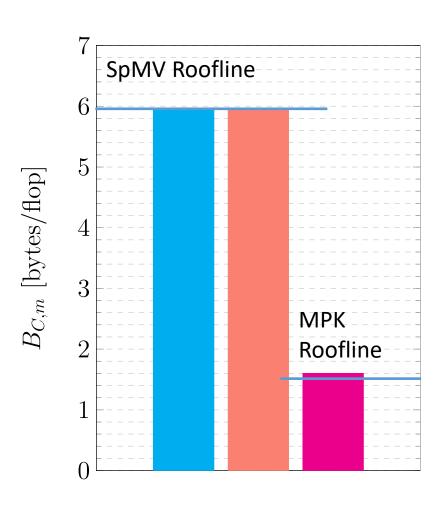
- Huber et al.: Graph-based higher-order time integration of PDEs¹
 - "Geometrical approach" based on matrix bandwidth
 - Works for 2D stencil matrices → Runs into problem for 3D and/or unstructured matrices
- Mohiyuddin et al.: Minimizing communication in sparse matrix solvers²
 - "Domain decomposition" of underlying graph
 - Requires "ghosting" → Indirect accesses or redundant copies of the matrix entries → Scalability!!
- → Exploit level structure in RACE for cache blocking!



¹Huber et al., 2021. Graph-based multi-core higher-order time integration of linear autonomous partial differential equations. J. Comput. Sci. <u>DOI:10.1016/j.jocs.2021.101349</u>
²Mohiyuddin et al., 2009. Minimizing communication in sparse matrix solvers. In Proceedings of the SC′209. <u>DOI:10.1145/1654059.1654096</u>

RACE MPK – First Implementation





Intel Xeon Gold 6248

1 Socket (20c)

pwtk matrix

- $N_r = 217,918$
- N_{nz} = 11,634,424

Performance

Memory traffic

RACE MPK - Performance Problem Identified

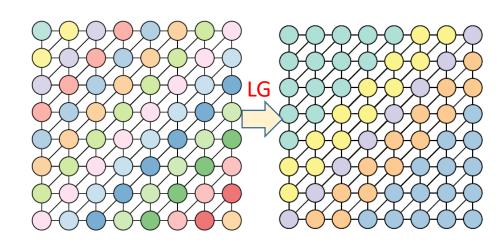
- Scheme seems to work (reduces data traffic) at least for pwtk
- But: Performance ⊗ !!!!

- Analysis of hardware performance counters (LIKWID) for pwtk matrix: INSTR RETIRED ANY up 2x for level based SpMV!
 - → Frequent thread syncronisations!

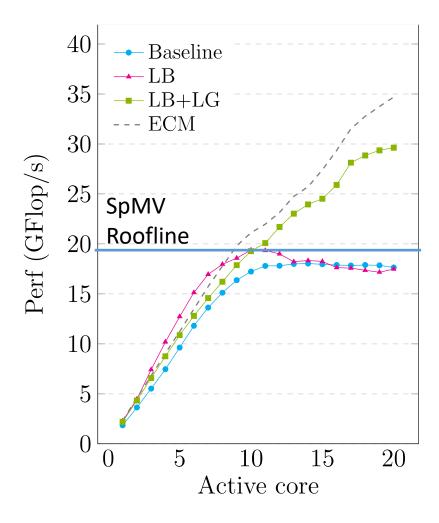
Reason: After each level threads sync!

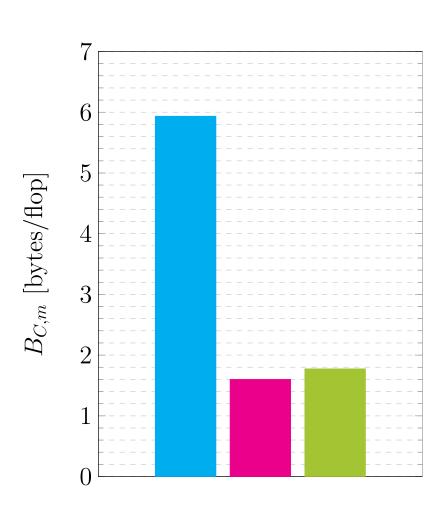
Measures:

- → Reduce #levels by level aggregation ("LG")
- → Global sync. replaced by point-to-point sync. ("p2p")



RACE MPK – LG optimization





Intel Xeon Gold 6248

• 1 Socket (20c)

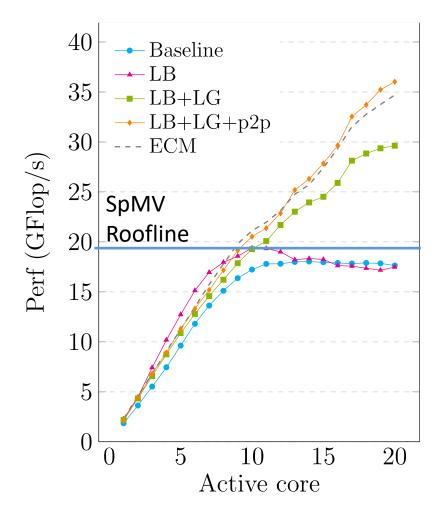
pwtk matrix

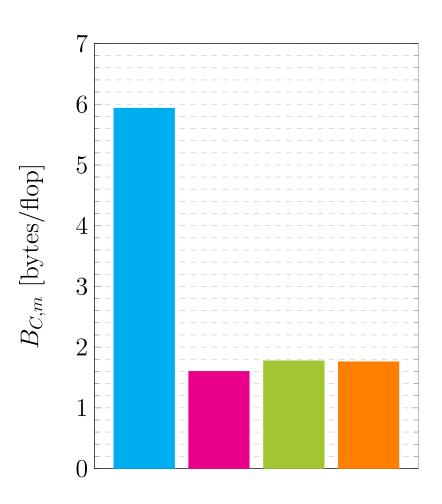
- $N_r = 217,918$
- $N_{nz} = 11,634,424$

Performance

Memory traffic

RACE MPK – LG+p2p optimization





Intel Xeon Gold 6248

• 1 Socket (20c)

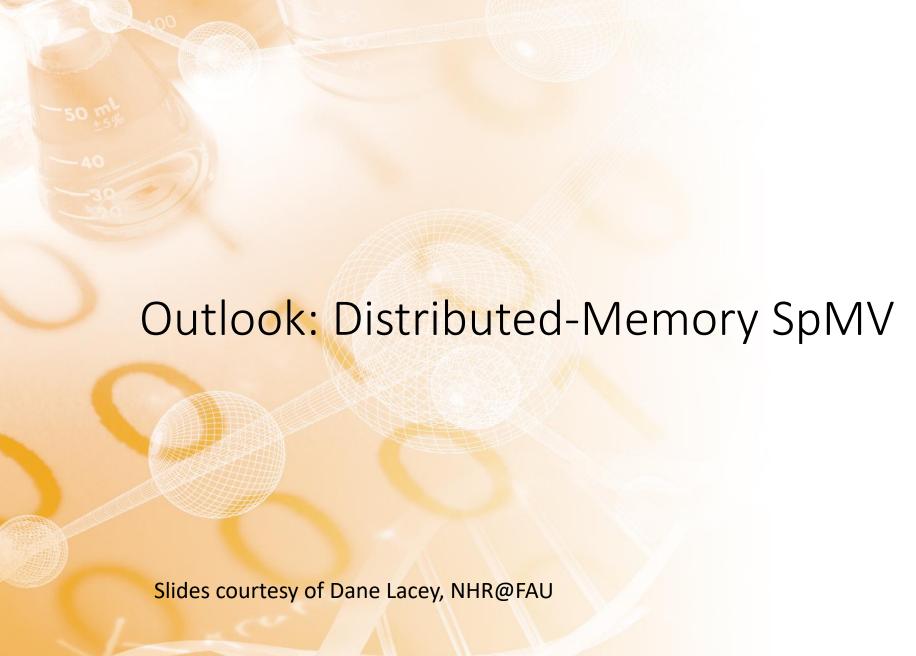
pwtk matrix

- $N_r = 217,918$
- N_{nz} = 11,634,424



Performance

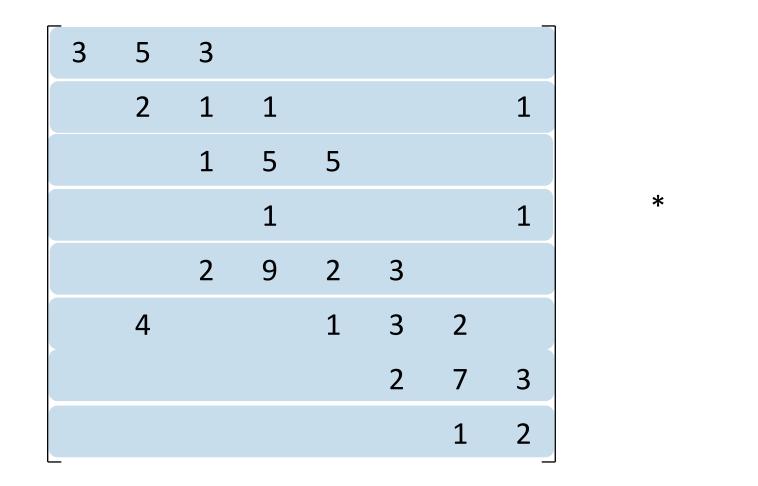
Memory traffic

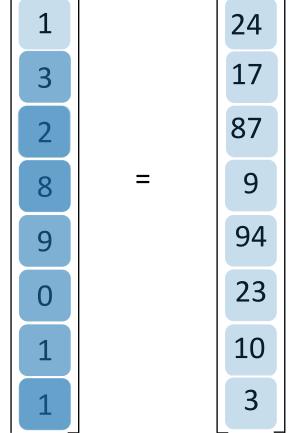


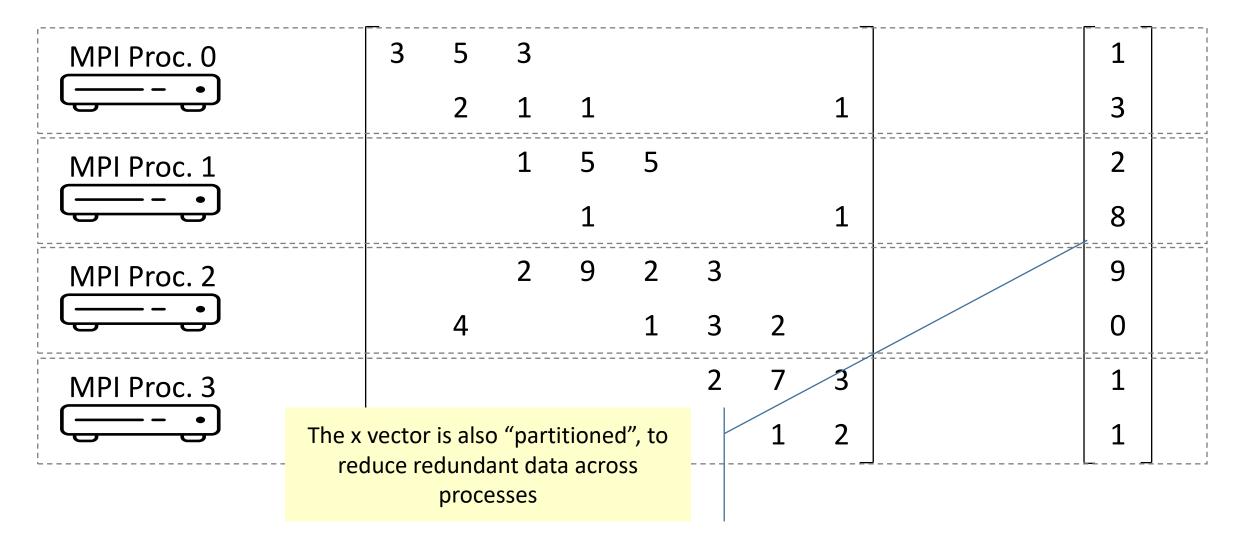
SpMV Example

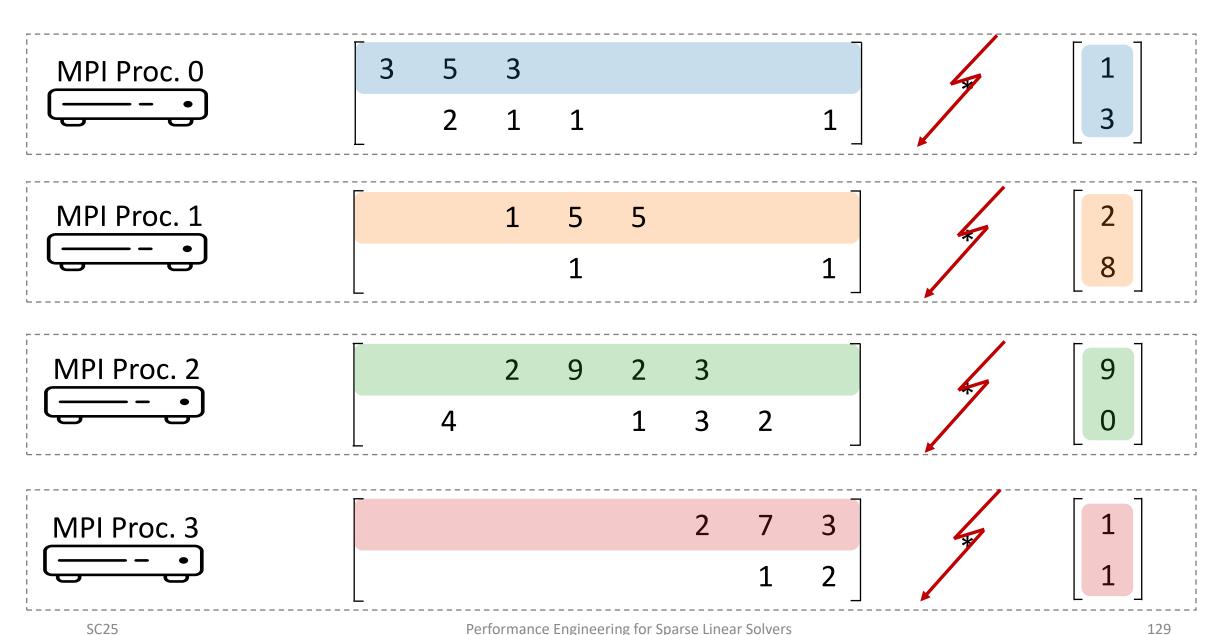
	3	5	3							1		0	
A =		2	1	1				1		3	y =	0	
		4	1	5	5					2		0	
				1		2 3		1	x =	8		0	
			2	9	2				^	9	,	0	
					1	3	2			0		0	
						2	7	3		1		0	
	_						1	2 _				0	

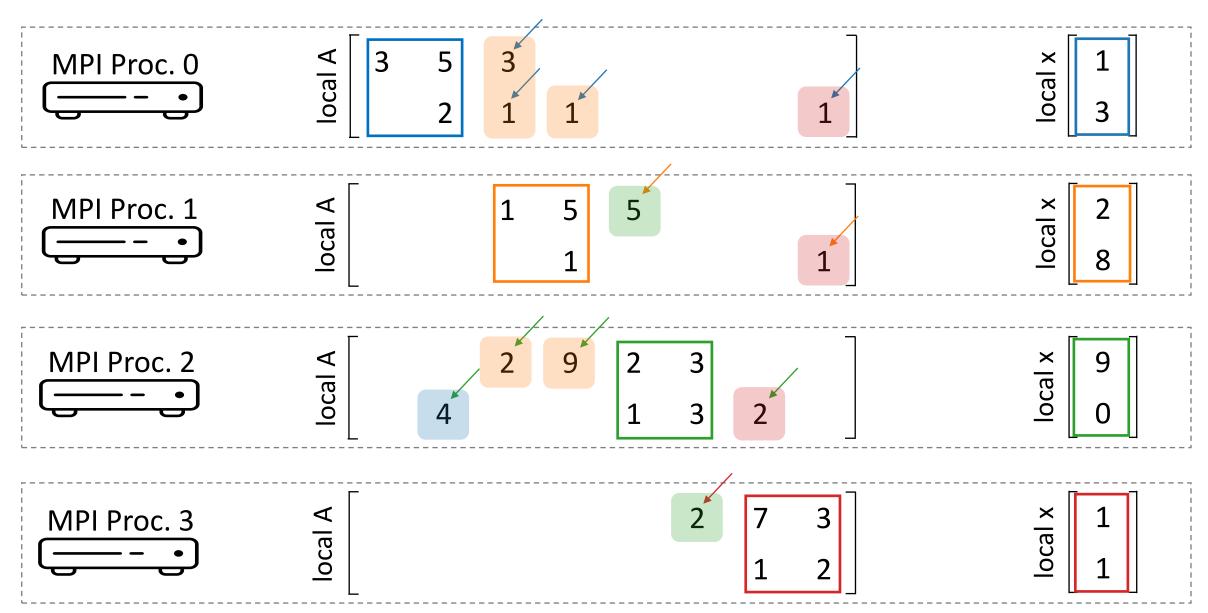
SpMV Example

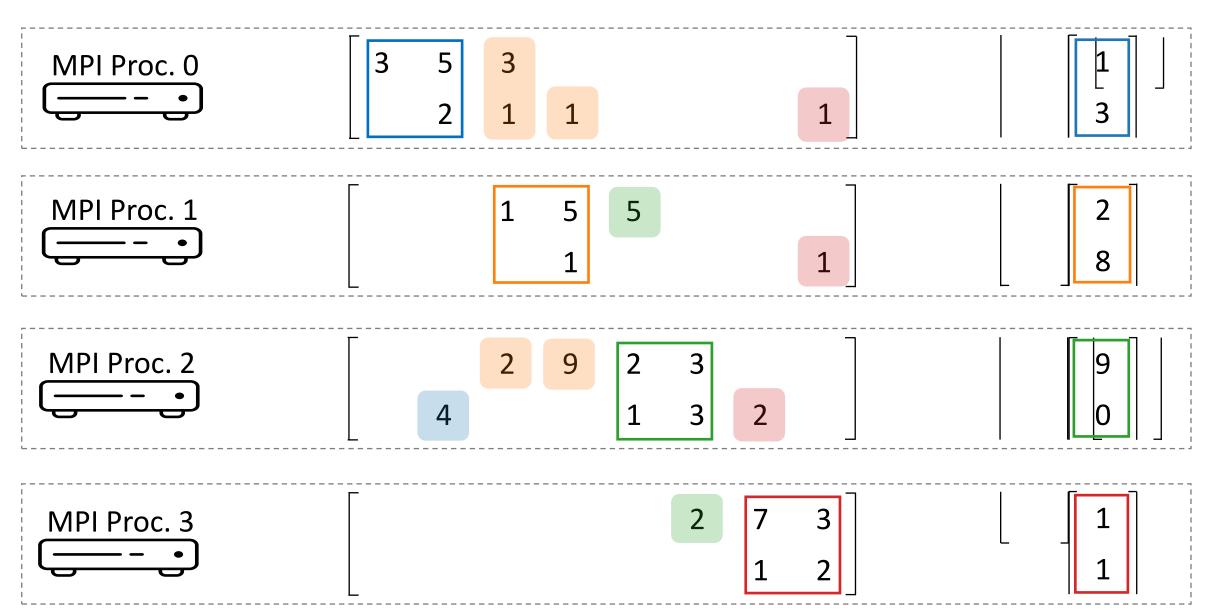


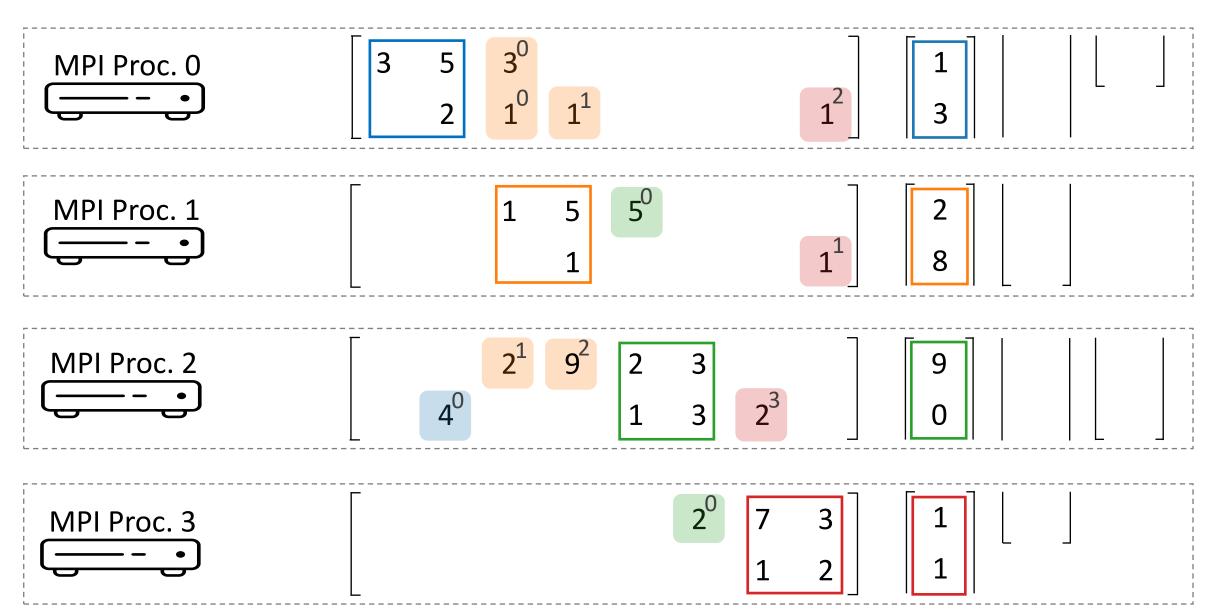


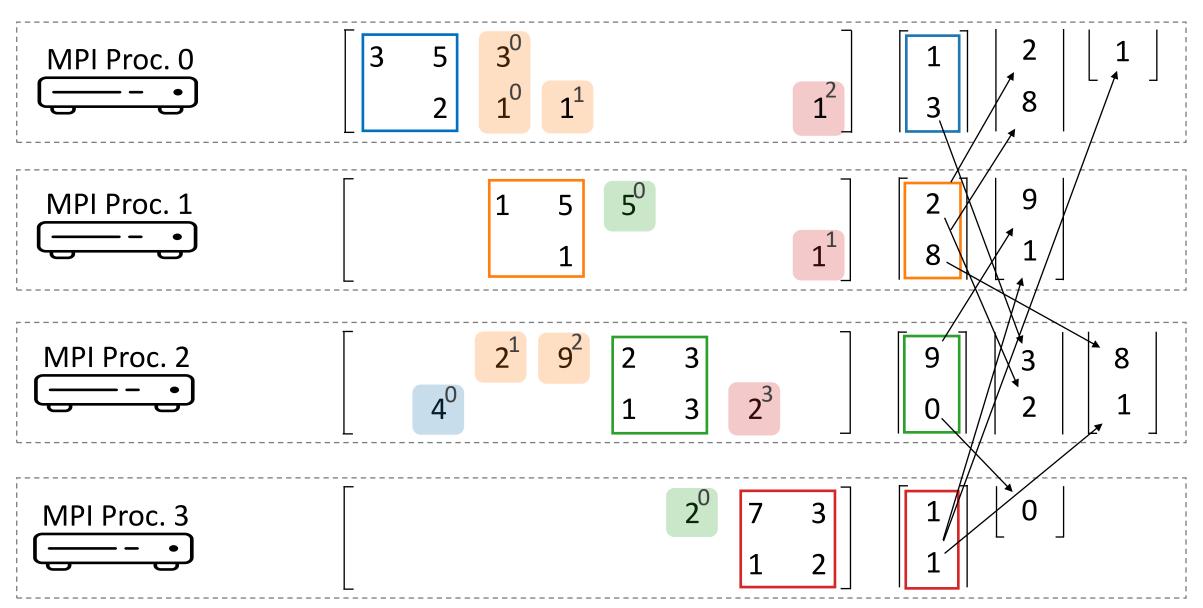


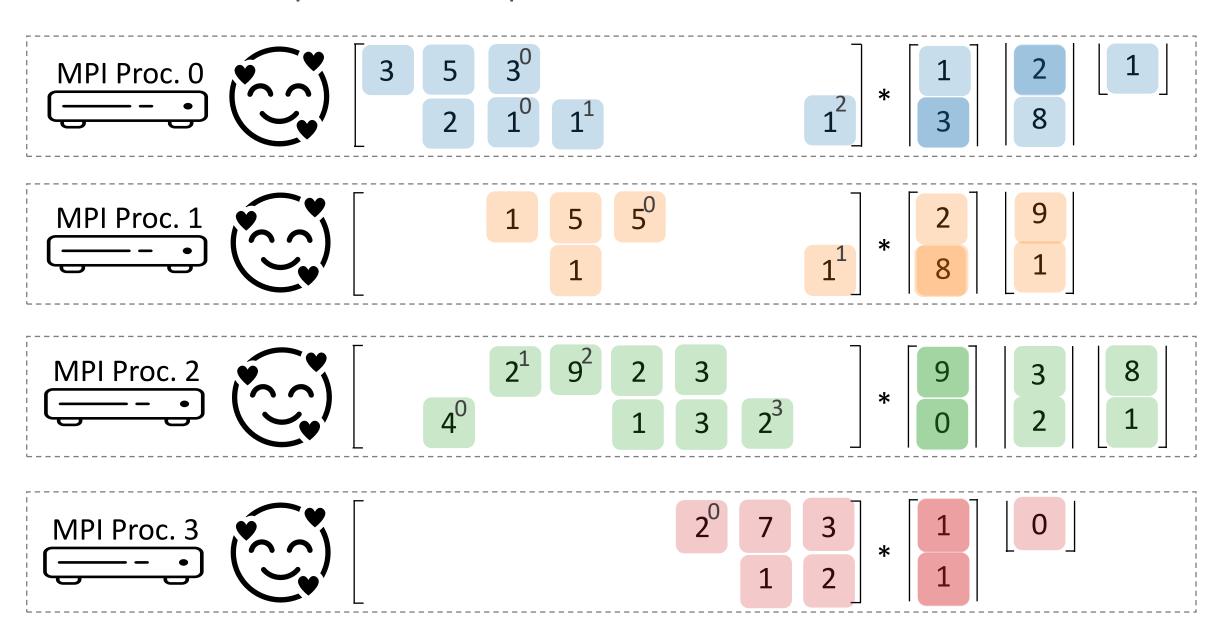


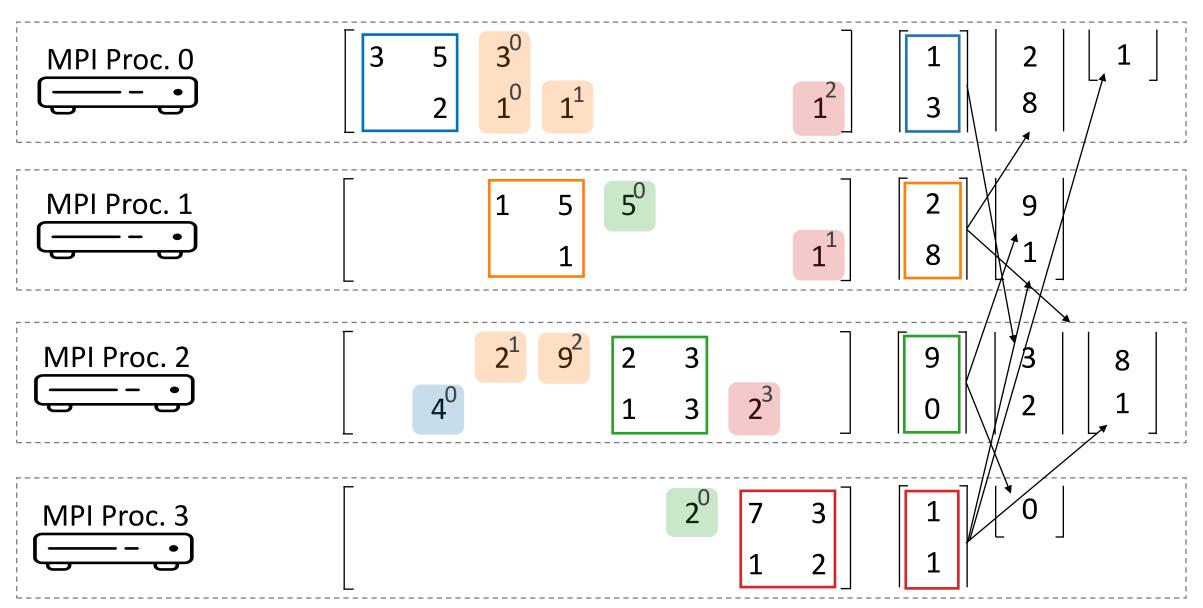


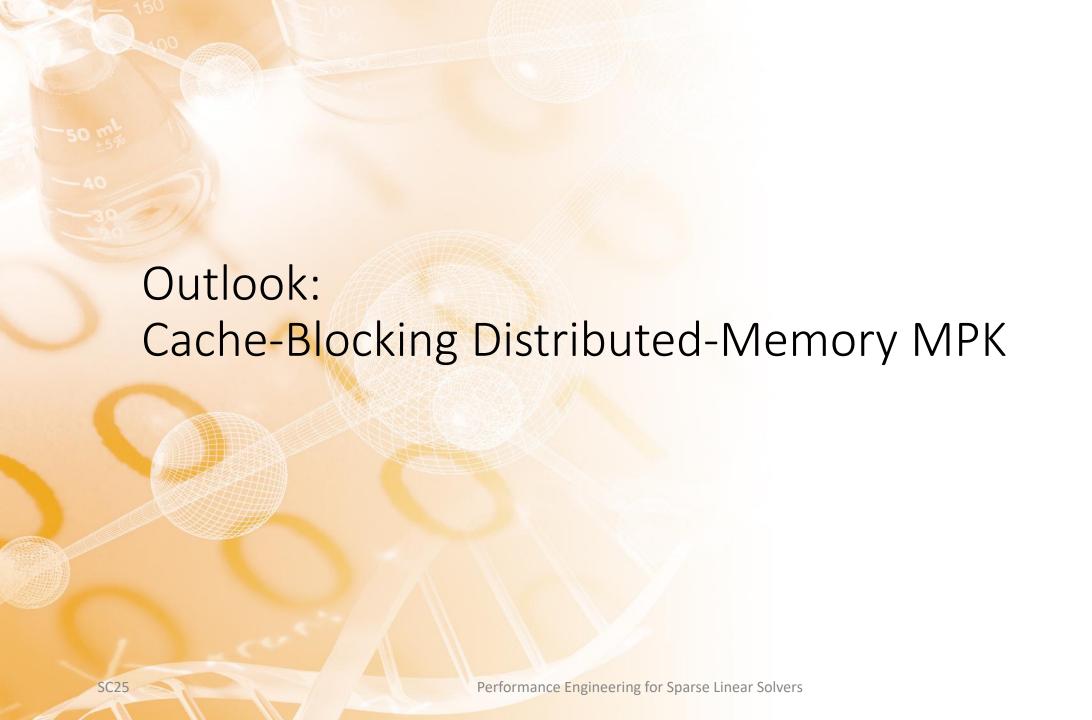


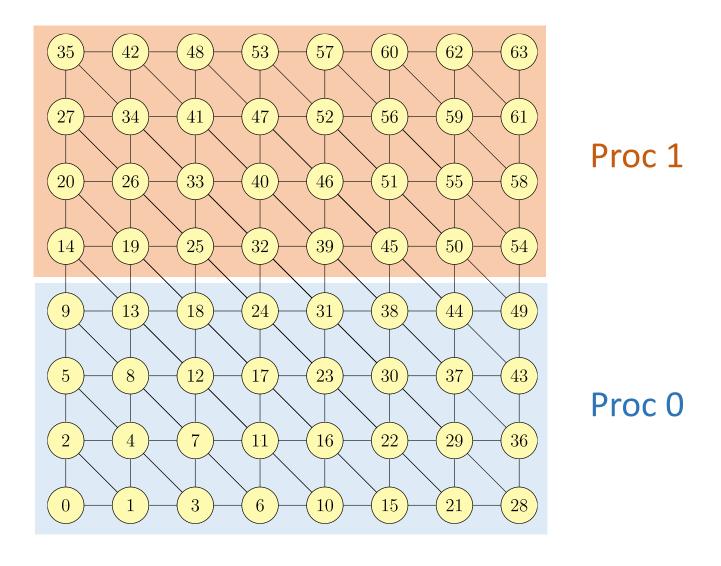


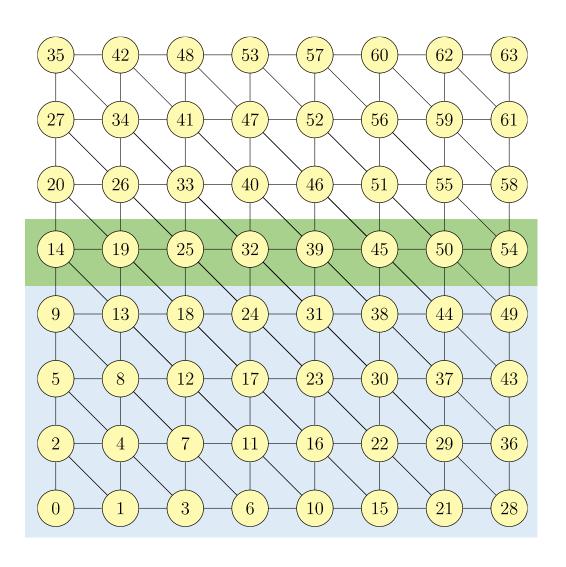






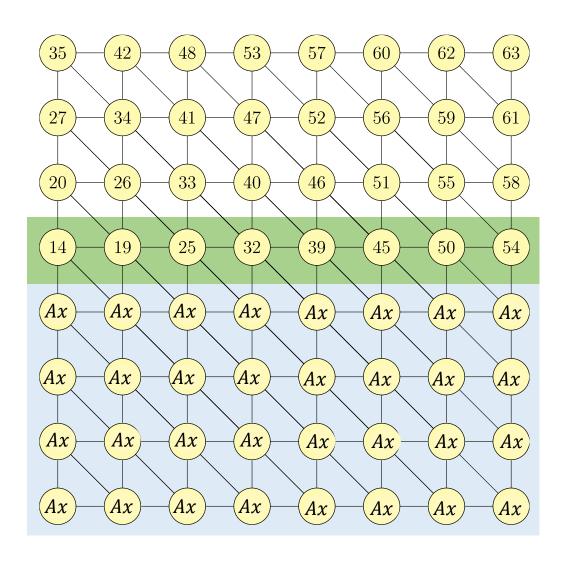


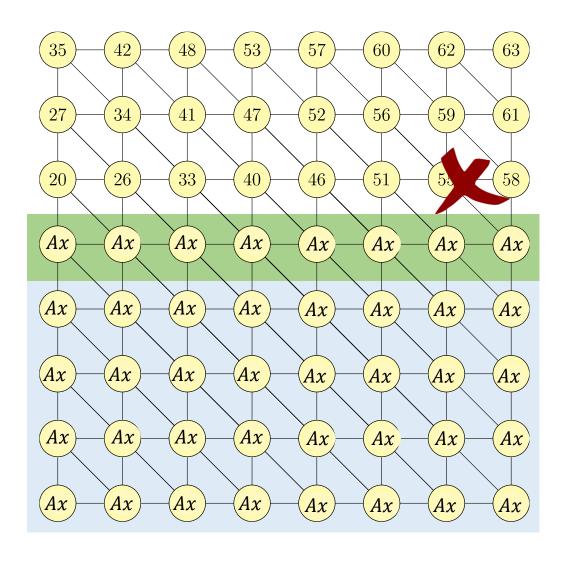




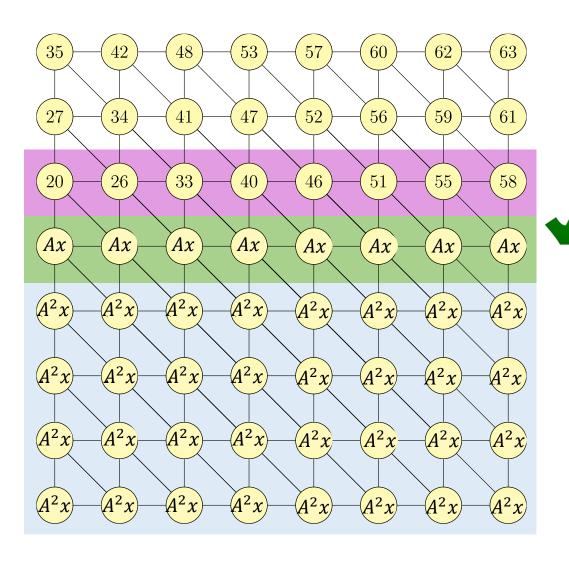
Computing Ax on Proc 0 requires neighbors of Proc 0.

How about computing $A^p x$?

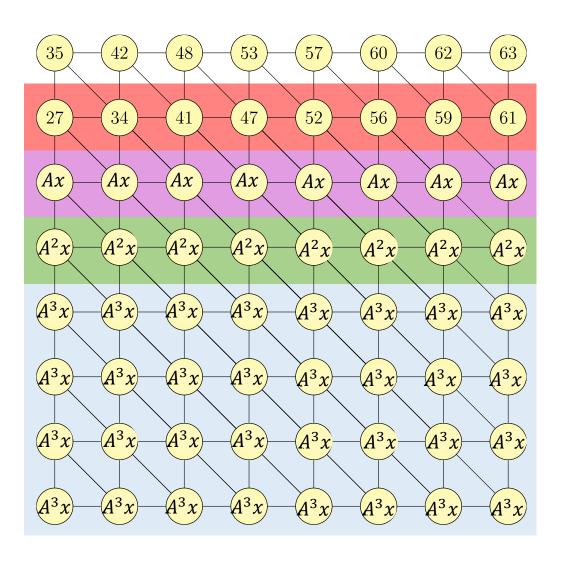




To compute A^2x we need Ax on the neighbors.

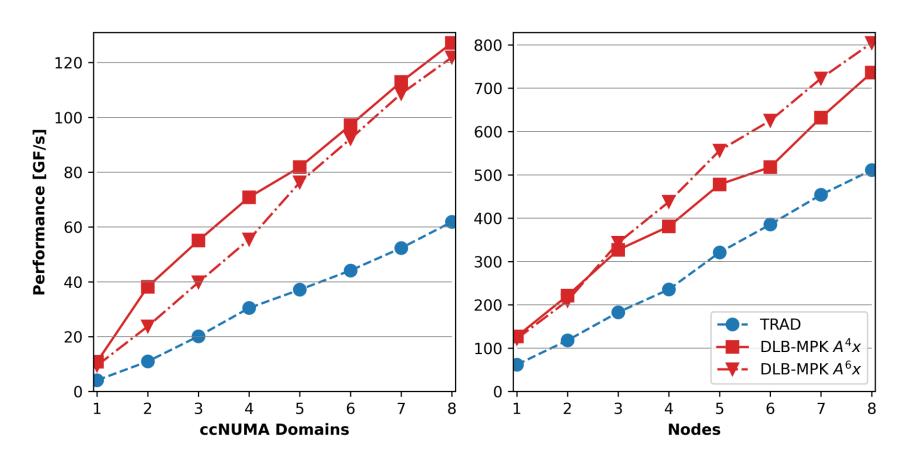


To compute A^2x we need Ax on the neighbors.



To compute A^3x we need A^2x on the neighbors.

In general, to compute $A^{p}x$ we need p neighbors.



It works!

No redundant work and/or extra communication required, see upcoming paper.*

^{*} D. Lacey, C. Alappat, F. Lange, G. Hager, and G. Wellein: *Cache Blocking of Distributed-Memory Parallel Matrix Power Kernels*, to be submitted.

Tutorial conclusions

- Memory bandwidth limitations are ubiquitous in sparse linear solvers
- SpMV performance depends on the storage format
- Roofline is an indispensable tool for performance analysis
- Time to solution is a fusion of flop/s performance and fast convergence
- Matrix powers can be optimized for better cache reuse



Performance Engineering for Linear Solvers

This tutorial covers code analysis, performance modeling, and optimization for linear solvers on CPU and GPU nodes. Performance Engineering is often taught using simple loops as instructive examples for performance models and how they can guide optimization; however, full, preconditioned linear solvers comprise multiple back-to-back loops enclosed in an iteration scheme that is executed until convergence is achieved. Consequently, the concept of "optimal performance" has to account for both hardware resource efficiency and iterative solver convergence. We convey a performance engineering process that is geared towards linear iterative solvers. After introducing basic notions of hardware organization and storage for dense and sparse data structures, we show how the Roofline performance model can be applied to such solvers in predictive and diagnostic ways and how it can be used to assess the hardware efficiency of a solver, covering important corner cases such as pure memory boundedness. Then we advance to the structure of preconditioned solvers, using the Conjugate Gradient Method (CG) algorithm as a leading example. Hotspots and bottlenecks of the complete solver are identified followed by the introduction of advanced performance optimization techniques like preconditioning and cache blocking.





Christie Alappat received his master's degree with honors from the Bavarian Graduate School of Computational Engineering, Friedrich-Alexander-Universität Erlangen-Nürnberg. He is in the final stages of completing his doctoral studies under the guidance of Prof. Gerhard Wellein. At the same time, he is currently working at Intel as a math algorithm engineer. His research interests include performance engineering, sparse matrix and graph algorithms, iterative linear solvers, and eigenvalue computation. He is the author of the RACE open-source software framework, which is used to accelerate challenging computations in sparse linear algebra on modern compute devices. He is also the lead author of a paper that received the SIAM Activity Group on Supercomputing (SIAG/SC) Best Paper Prize in 2024.

https://hpc.fau.de/person/christie-alappat/





Jonas has more than 20 years of experience in HPC and scientific computing with applications in CFD, climate research and quantum physics. Specifically, he has worked on domain decomposition methods for sparse linear systems, implicit ocean models, sparse eigenvalue problems on heterogeneous supercomputers, code optimization for multi-core CPUs and vector processors, and software and performance engineering for scientific applications.

Jonas has a PhD in applied mathematics (Groningen 2011). He spent two years at the Center for Interdisciplinary Mathematics in Uppsala, after which he moved to Cologne as a Scientific Employee of the German Aerospace Center (DLR) Institute for Software Technology. There he led a research group on parallel numerics from 2017 to 2021. Since June 2021 he is an Assistant Professor at the Delft High Performance Computing Center DHPC, where he coordinates the center's training activities.

https://www.tudelft.nl/en/eemcs/the-faculty/departments/applied-mathematics/people/dr-j-jonas-thies





Hartwig Anzt is the Chair of Computational Mathematics at the TUM School of Computation, Information and Technology of the Technical University of Munich (TUM) Campus Heilbronn. He also holds a Research Associate Professor position at the Innovative Computing Lab (ICL) at the University of Tennessee (UTK). Hartwig Anzt received a PhD in applied mathematics from the Karlsruhe Institute of Technology (KIT) and specializes in iterative methods and preconditioning techniques for the next generation hardware architectures. He also has a long track record of high-quality development. He is author of the MAGMA-sparse open-source software package and managing lead of the Ginkgo math software library. Hartwig Anzt had served as a PI in the Software Technology (ST) pillar of the US Exascale Computing Project (ECP), including a coordinated effort aiming at integrating low-precision functionality into high-accuracy simulation codes. He also is a PI in the EuroHPC project MICROCARD.

Hartwig Anzt is the main author of more than 100 peer-reviewed publications, part of the scientific committee of international conferences, Associate Editor of the SIAM Journal on Scientific Computing (SISC), Associate Editor of ACM Transactions on Parallel Computing, workshop chair for ISC High Performance 2022, and has been elected as SIAM Activity Group on Supercomputing program manager.

https://hartwiganzt.github.io/





Georg Hager holds a PhD and a habilitation degree in Computational Physics from the University of Greifswald. Since 2021 he heads the Training and Support Division of the newly founded "Erlangen National High Performance Computing Center" (NHR@FAU). Previously he was a senior researcher in the HPC Services group at Erlangen Regional Computing Center (RRZE), which is part of the Friedrich-Alexander-Universität Erlangen-Nürnberg. Recent research includes architecture-specific optimization strategies for current microprocessors, performance engineering of scientific codes, and analytic modeling of massively parallel programs. His textbook "Introduction to High Performance Computing for Scientists and Engineers" is recommended or required reading in many HPC-related lectures and courses worldwide. He has more than two decades of experience in teaching high performance computing and performance engineering to students and scientists. Together with colleagues from NHR@FAU and other centers, he conducts long-standing series of tutorials on Performance Engineering and Hybrid Programming.

https://blogs.fau.de/hager