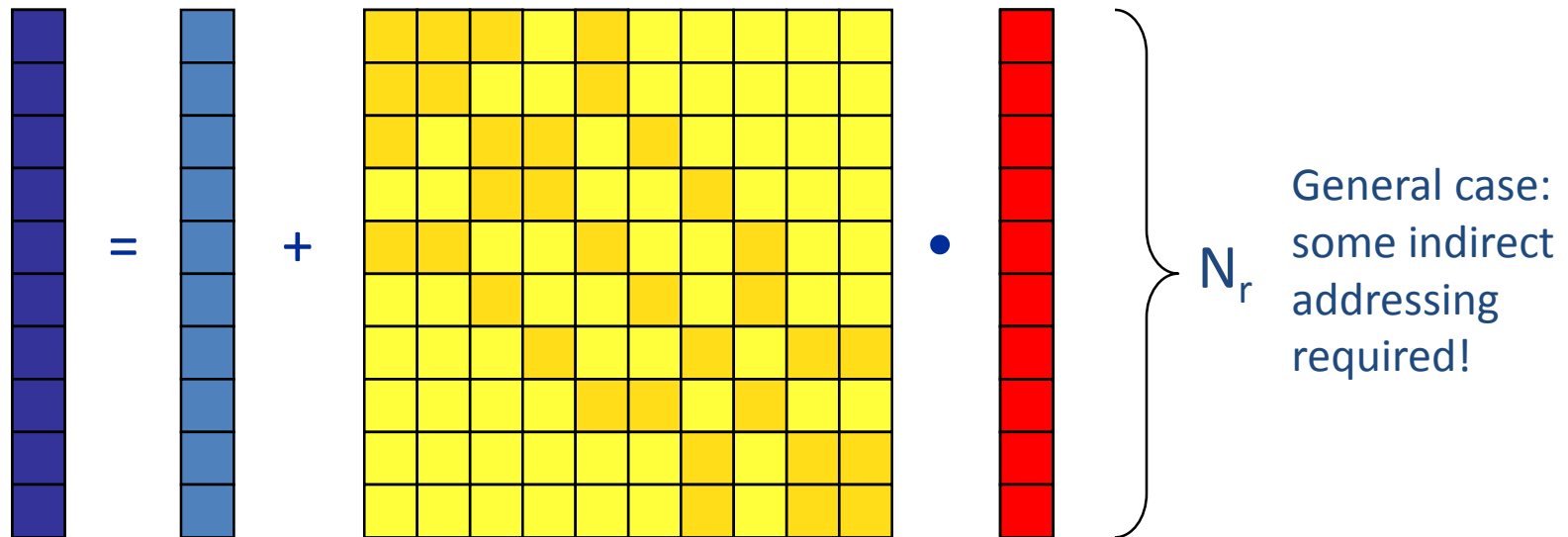
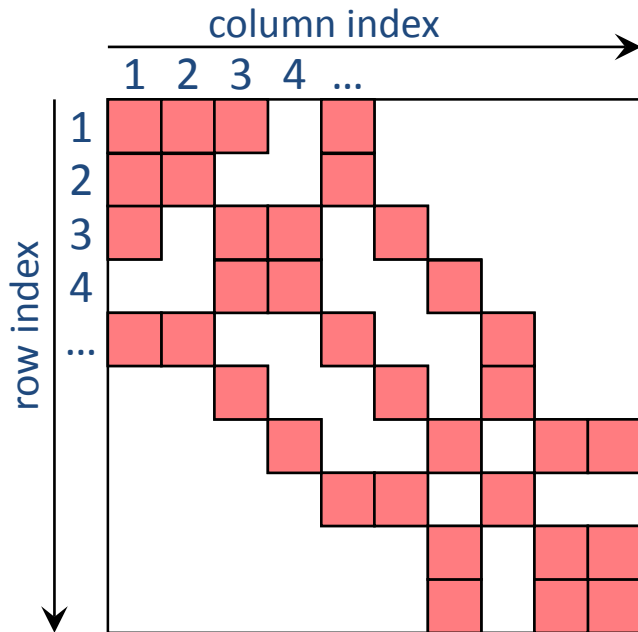


Case study: Sparse Matrix-Vector Multiplication

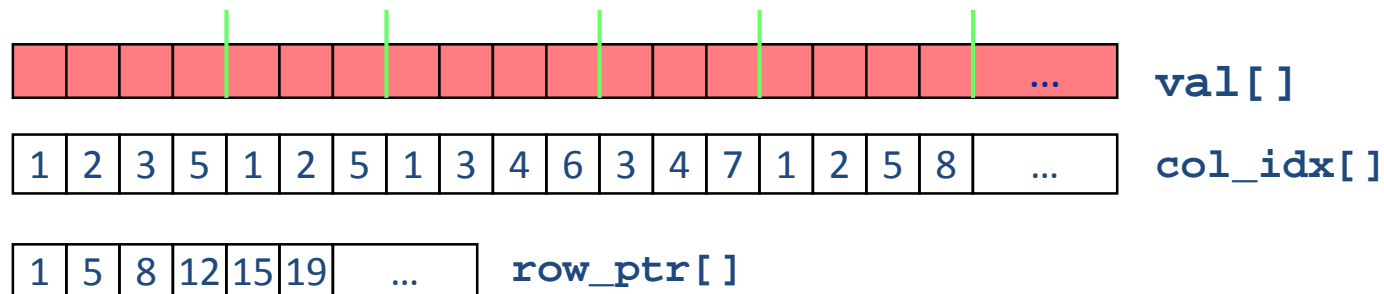
- Key ingredient in some matrix diagonalization algorithms
 - Lanczos, Davidson, Jacobi-Davidson
- Store only N_{nz} nonzero elements of matrix and RHS, LHS vectors with N_r (number of matrix rows) entries
- “Sparse”: $N_{nz} \sim N_r$
- Average number of nonzeros per row: $N_{nzs} = N_{nz}/N_r$



- For large problems, SpMV is inevitably **memory-bound**
 - **Intra-socket saturation effect** on modern multicores
- SpMV is **easily parallelizable** in shared and distributed memory
 - Load balancing
 - Communication overhead
- Data storage format is **crucial** for performance properties
 - Most useful general format on CPUs:
Compressed Row Storage (**CRS**)
 - Depending on compute architecture



- **val[]** stores all the nonzeros (length N_{nz})
- **col_idx[]** stores the column index of each nonzero (length N_{nz})
- **row_ptr[]** stores the starting index of each new row in **val[]** (length: N_r)



- Strongly memory-bound for large data sets
 - Streaming, with partially indirect access:

```
!$OMP parallel do schedule(???)  
do i = 1, Nr  
  do j = row_ptr(i), row_ptr(i+1) - 1  
    C(i) = C(i) + val(j) * B(col_idx(j))  
  enddo  
enddo  
!$OMP end parallel do
```

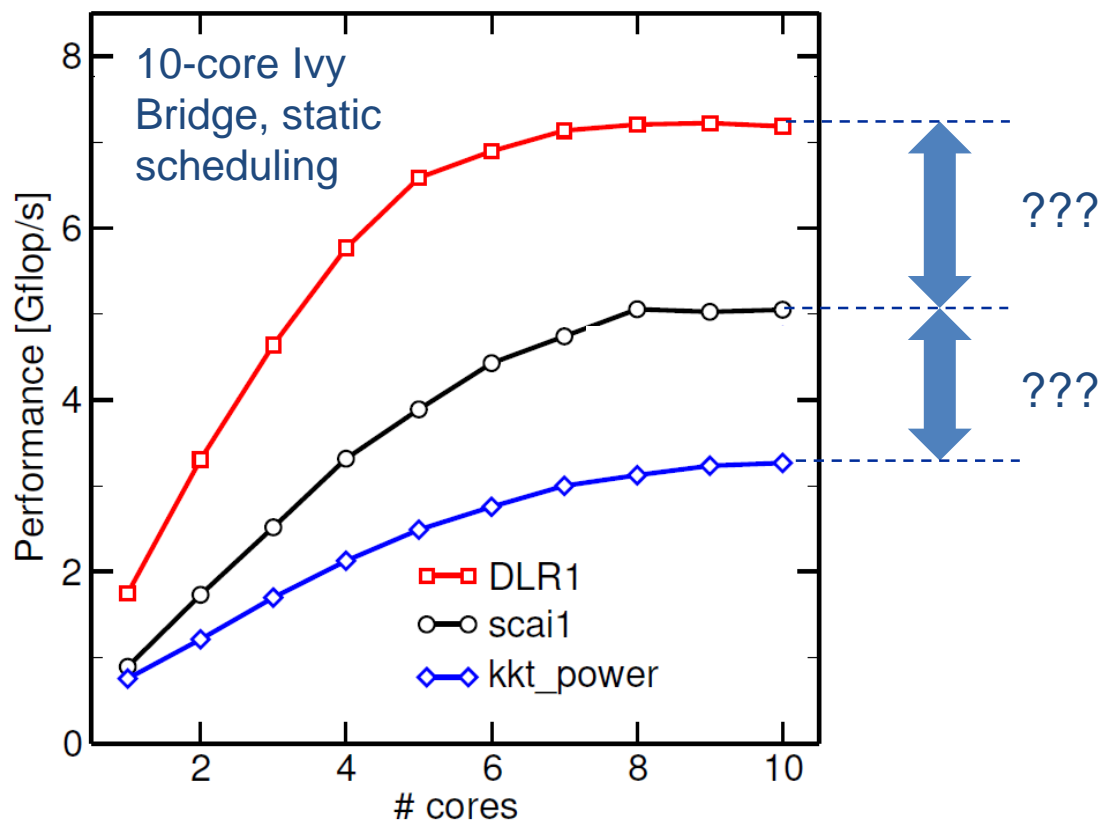
- Usually many spMVMs required to solve a problem
- Now let's look at some performance measurements...

- Strongly memory-bound for large data sets → **saturation performance** across cores on the chip
- Performance seems to depend on the matrix

- Can we explain this?

- Is there a “light speed” for SpMV?

- Optimization?



```
do i = 1, Nr
  do j = row_ptr(i), row_ptr(i+1) - 1
    C(i) = C(i) + val(j) * B(col_idx(j))
  enddo
enddo
```

```
real*8      val[Nnz]
integer*4   col_idx[Nnz]
integer*4   row_ptr[Nr]
real*8      C[Nr]
real*8      B[Nc]
```

Min. load traffic [B]: $(8 + 4) N_{nz} + (4 + 8) N_r + 8 N_c$

Min. store traffic [B]: $8 N_r$

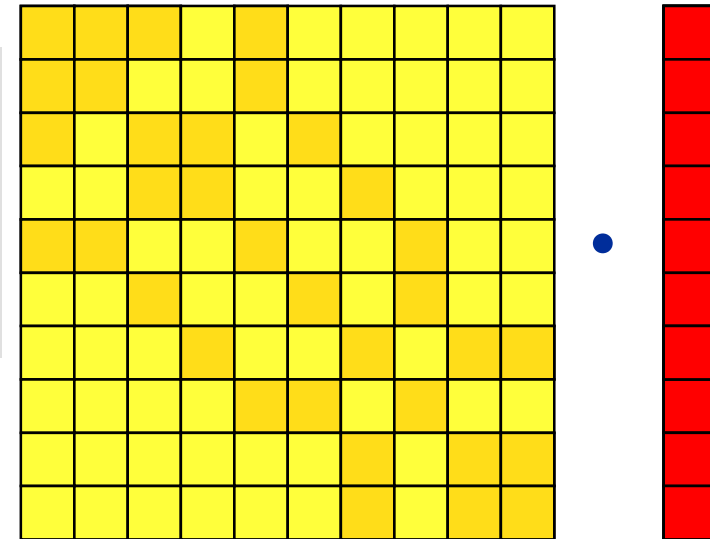
Total FLOP count [F]: $2 N_{nz}$

$$B_{C,min} = \frac{12 N_{nz} + 20 N_r + 8 N_c}{2 N_{nz}} \frac{B}{F} = \frac{12 + 20/N_{nzc} + 8/N_{nzc}}{2} \frac{B}{F}$$

Nonzeros per row ($N_{nzc} = N_{nz}/N_c$) or column ($N_{nzc} = N_{nz}/N_c$)

Lower bound for code balance: $B_{C,min} \geq 6 \frac{B}{F} \rightarrow I_{max} \leq \frac{1}{6} \frac{F}{B}$

```
do i = 1, Nr
  do j = row_ptr(i), row_ptr(i+1) - 1
    C(i) = C(i) + val(j) * B(col_idx(j))
  enddo
enddo
```



$$B_{C,min} = \frac{12 + 20/N_{nzc} + 8/N_{nzc}}{2} \frac{B}{F}$$

$$B_C(\alpha) = \frac{12 + 20/N_{nzc} + 8\alpha}{2} \frac{B}{F}$$

Parameter (α) quantifies additional traffic for $B(\cdot)$ (irregular access):

$$\alpha \geq 1/N_{nzc}$$

$$\alpha N_{nzc} \geq 1$$

Consider square matrices: $N_{nzc} = N_{nzc}$ and $N_c = N_r$

Note: $B_C(1/N_{nzc}) = B_{C,min}$

DP CRS code balance

- α quantifies the traffic for loading the RHS
 - $\alpha = 0$ → RHS is in cache
 - $\alpha = 1/N_{nzs}$ → RHS loaded once
 - $\alpha = 1$ → no cache
 - $\alpha > 1$ → Houston, we have a problem!
- “Target” performance = b_S/B_C
- **Caveat:** Maximum memory BW may not be achieved with spMVM (see later)

$$B_C(\alpha) = \frac{12 + 20/N_{nzs} + 8\alpha B}{2} \frac{B}{F}$$
$$= \left(6 + 4\alpha + \frac{10}{N_{nzs}}\right) \frac{B}{F}$$

Can we predict α ?

- Not in general
- Simple cases (banded, block-structured): Similar to layer condition analysis

→ Determine α by measuring **the actual memory traffic**
(→ measured code balance B_C^{meas})

$$B_C(\alpha) = \left(6 + 4\alpha + \frac{10}{N_{nzs}}\right) \frac{B}{F} = \frac{V_{meas}}{N_{nzs} \cdot 2 F} (= B_C^{meas})$$

- V_{meas} is the measured overall memory data traffic (using, e.g., likwid-perfctr)

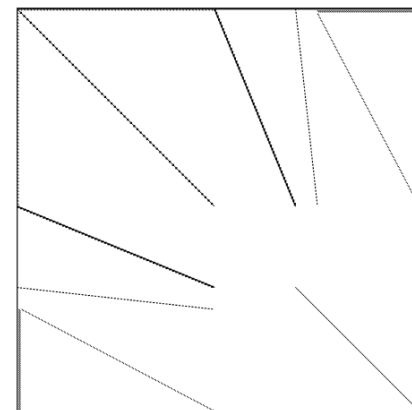
- Solve for α :
$$\alpha = \frac{1}{4} \left(\frac{V_{meas}}{N_{nzs} \cdot 2 \text{ bytes}} - 6 - \frac{10}{N_{nzs}} \right)$$

Example: kkt_power matrix from the UoF collection on one Intel SNB socket

- $N_{nzs} = 14.6 \cdot 10^6, N_{nzs} = 7.1$
- $V_{meas} \approx 258 \text{ MB}$
- $\rightarrow \alpha = 0.36, \alpha N_{nzs} = 2.5$
- \rightarrow RHS is loaded 2.5 times from memory
- and:

$$\frac{B_C(\alpha)}{B_{C,min}} = 1.11$$

11% extra traffic \rightarrow
optimization potential!

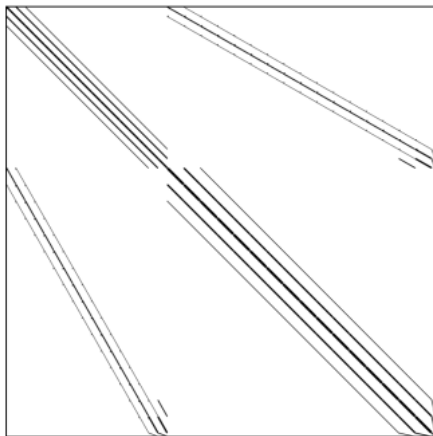


Three different sparse matrices

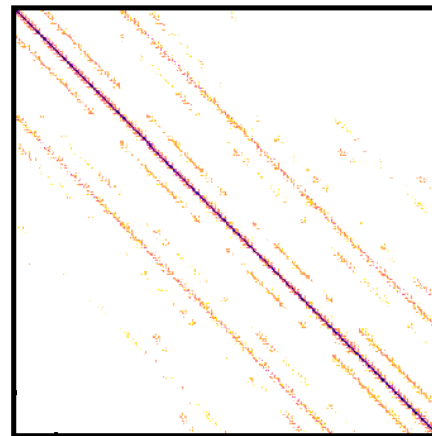
Benchmark system: Intel Xeon Ivy Bridge E5-2660v2, 2.2 GHz, $b_S = 46.6$ GB/s

$$\rightarrow \text{Roofline: } P_{opt} = \frac{b_S}{B_{C,min}}$$

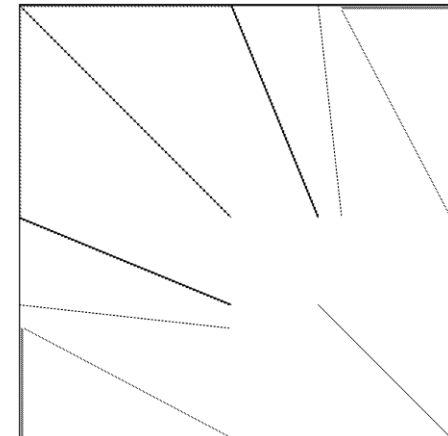
Matrix	N	N_{nzs}	$B_{C,min}$ [B/F]	P_{opt} [GF/s]
DLR1	278,502	143	6.1	7.64
scai1	3,405,035	7.0	8.0	5.83
kkt_power	2,063,494	7.08	8.0	5.83



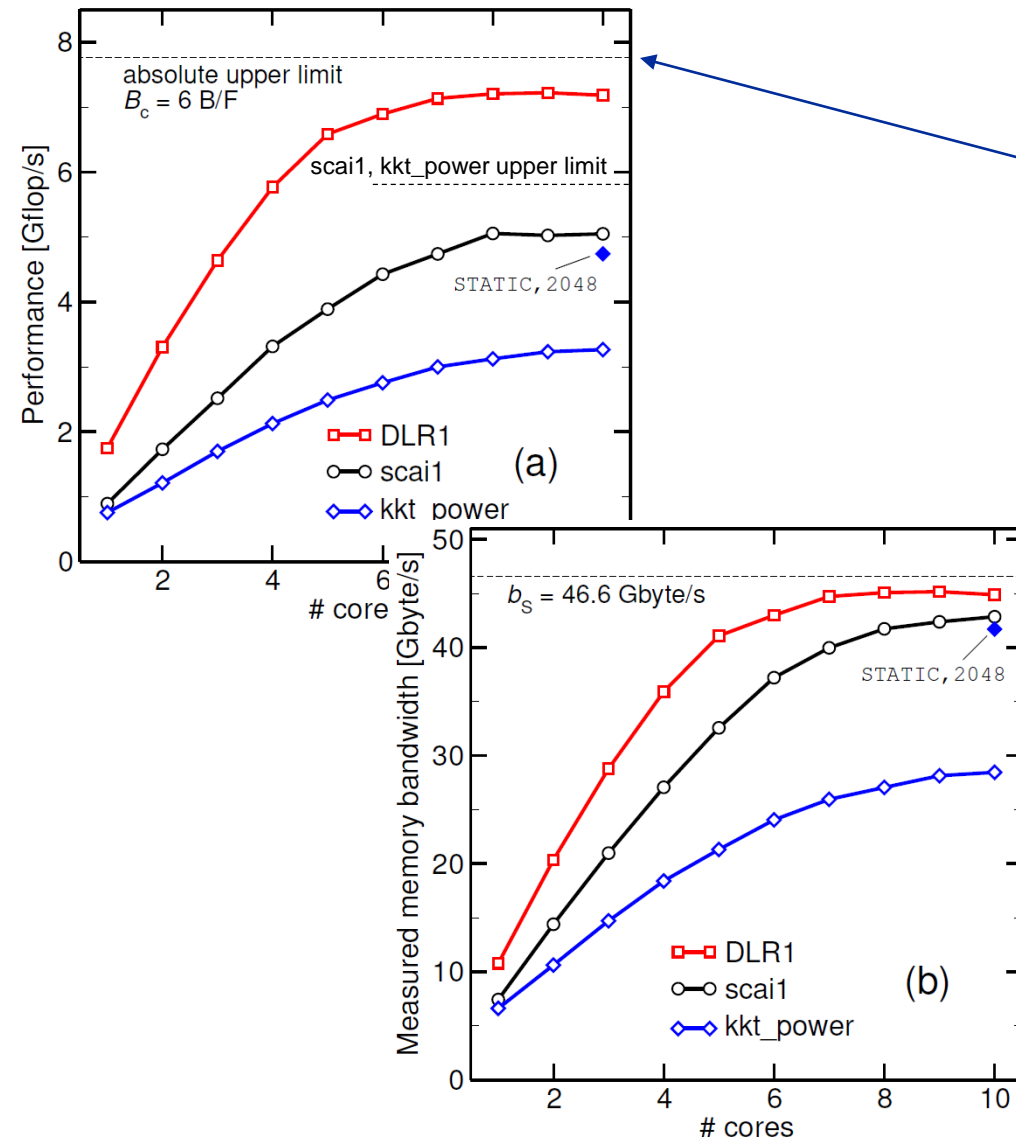
DLR1



scai1



kkt_power



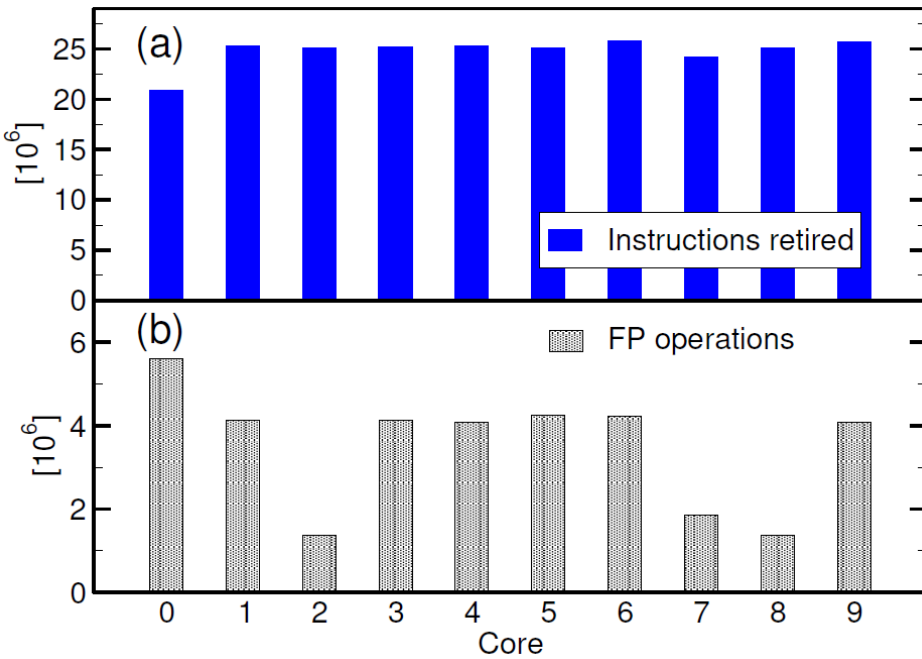
- $b_S = 46.6 \text{ GB/s}$, $B_c = 6 \text{ B/F}$
- Maximum spMVM performance:

$$P_{max} = 7.8 \text{ GF/s}$$

- **DLR1** causes minimum CRS code balance (as expected)
- **scai1** measured balance:

$$B_c^{meas} \approx 8.5 \text{ B/F} > B_{C,min}$$
 → good BW utilization, slightly non-optimal α
- **kkt_power** measured balance:

$$B_c^{meas} \approx 8.8 \text{ B/F} > B_{C,min}$$
 → performance degraded by load imbalance, fix by block-cyclic schedule

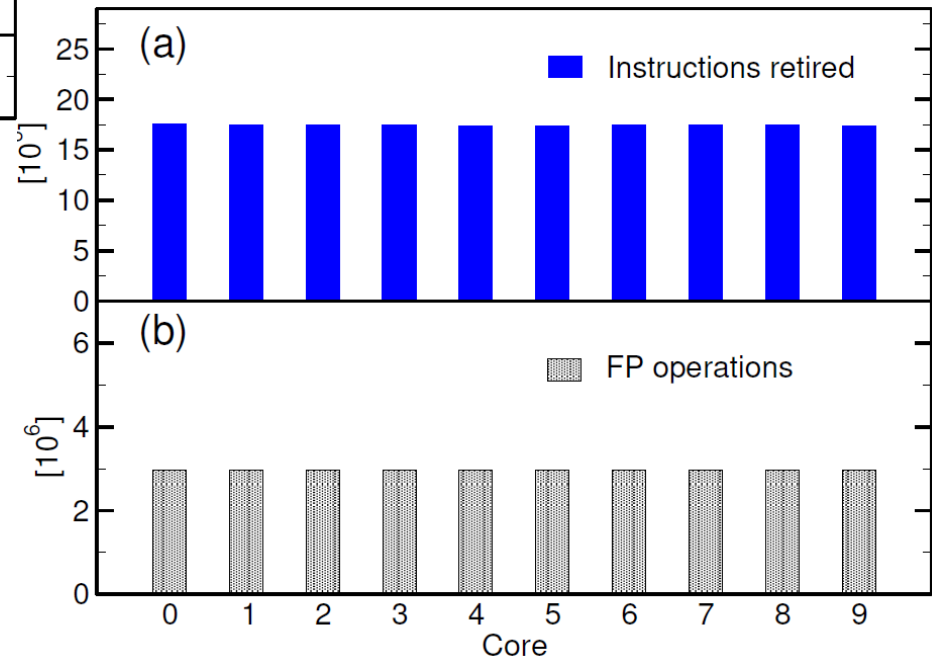


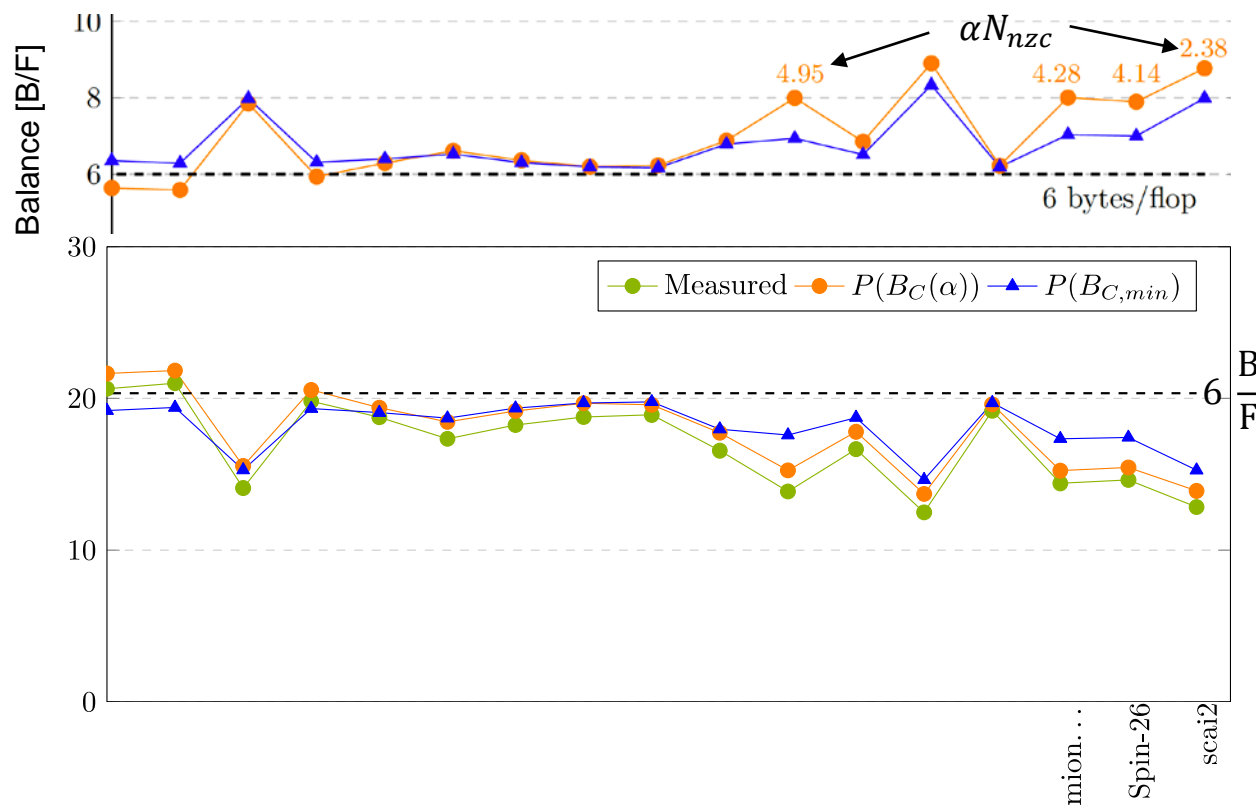
Measurements with likwid-perfctr (MEM_DP group)



→ Fewer overall instructions, (almost) BW saturation, 50% better performance with load balancing

→ **CPI value unchanged!**



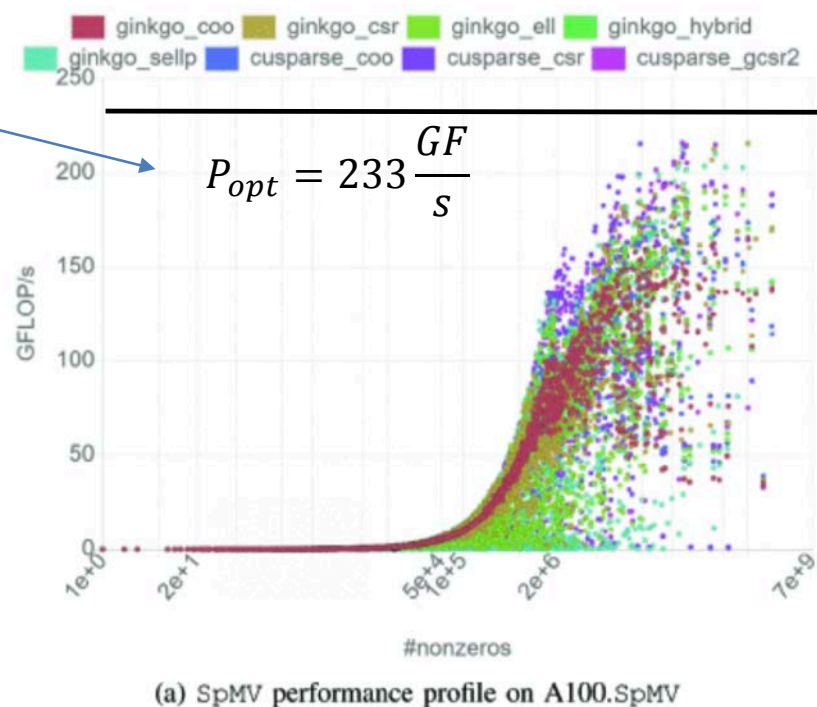
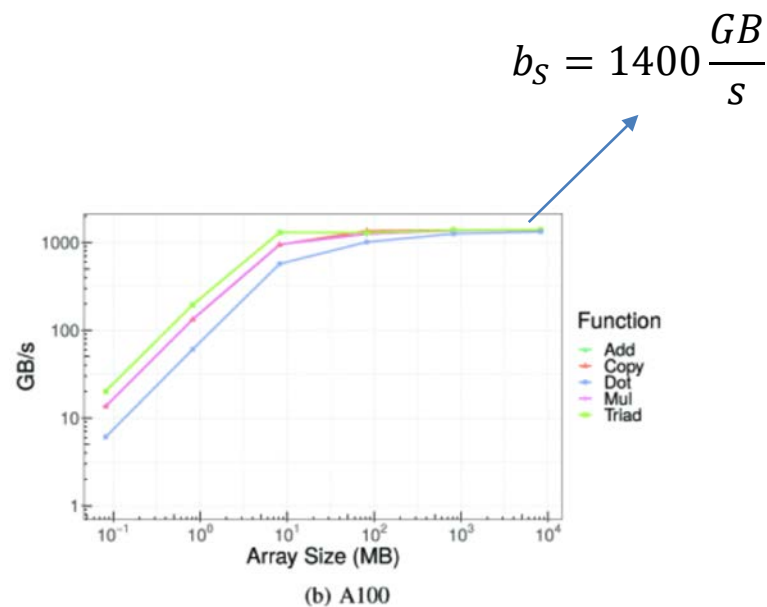


Intel Xeon Platinum 9242
24c@2.8GHz (turbo)
 $b_S = 122 \text{ GB/s}$

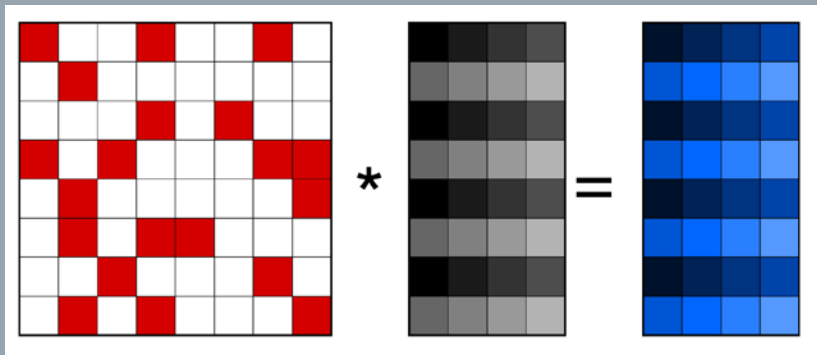
Matrices taken from: C. L. Alappatet al.: *ECM modeling and performance tuning of SpMV and Lattice QCD on A64FX*. Submitted. Preprint: [arXiv:2103.0301](https://arxiv.org/abs/2103.0301)

SpMV on NVIDIA A100:

- Different data formats and libraries
- 2800 of matrices (SuiteSparse Matrix collection)



Applying sparse matrix to multiple vectors (Sparse Matrix Multiple Vectors; SpMMV)



Multiple RHS vectors (SpMMV)

Unchanged matrix applied to multiple RHS (\mathbf{r}) vectors to yield multiple LHS (\mathbf{r}) vectors

```
do s = 1,r
  do i = 1, Nr
    do j = row_ptr(i),row_ptr(i+1)-1
      C(i,s) = C(i,s) + val(j) *
                B(col_idx(j),s)
    enddo
  enddo
enddo
```

B_c unchanged, no reuse of matrix data

```
do i = 1, Nr
  do j = row_ptr(i),row_ptr(i+1)-1
    do s = 1,r
      C(i,s) = C(i,s) + val(j) *
                B(col_idx(j),s)
    enddo
  enddo
enddo
```

Higher B_c due to max reuse of matrix data

```
do i = 1, Nr
  do j = row_ptr(i),row_ptr(i+1)-1
    do s = 1,r
      C(s,i) = C(s,i) + val(j) *
                B(s,col_idx(j))
    enddo
  enddo
enddo
```

CL-friendly data structure (row major)

One complete inner (**s**) loop traversal:

- $2r$ flops
- 12 bytes from matrix data (value + index)
- $\frac{16r}{N_{nzs}}$ bytes from the r LHS updates
- $\frac{4}{N_{nzs}}$ bytes from the row pointer
- $8r\alpha(r)$ bytes from the r RHS reads

```
do i = 1, Nr
  do j = row_ptr(i), row_ptr(i+1)-1
    do s = 1, r
      C(s,i) = C(s,i) + val(j) *
                B(s,col_idx(j))
    enddo
  enddo
enddo
```

$$B_c(r) = \frac{1}{2r} \left(12 + 8r\alpha(r) + \frac{16r + 4}{N_{nzs}} \right) \frac{B}{F}$$

$$= \left(\frac{6}{r} + 4\alpha(r) + \frac{8 + 2/r}{N_{nzs}} \right) \frac{B}{F}$$

OK so what now???

Let's check some limits to see if this makes sense!

$$B_c(r) = \left(\frac{6}{r} + 4\alpha(r) + \frac{8 + 2/r}{N_{nzs}} \right) \frac{B}{F}$$

$\xrightarrow{r=1}$ $\left(6 + 4\alpha + \frac{10}{N_{nzs}} \right) \frac{B}{F}$
reassuring 😊

$\xrightarrow{r \gg 1}$ $\left(4\alpha(r) + \frac{8}{N_{nzs}} \right) \frac{B}{F}$
Can become very small for large $N_{nzs} \rightarrow$ decoupling from memory bandwidth is possible!

$\xrightarrow{N_{nzs} \gg 1}$ $\frac{6}{r} \frac{B}{F}$

M. Kreutzer et al.: *Performance Engineering of the Kernel Polynomial Method on Large-Scale CPU-GPU Systems*. Proc. [IPDPS15](#), DOI: [10.1109/IPDPS.2015.76](#)

- **Conclusion from the Roofline analysis**
 - The roofline model does not “work” for spMVM due to the RHS traffic uncertainties
 - We have “**turned the model around**” and measured the actual memory traffic to determine the RHS overhead
 - Result indicates:
 1. how much actual traffic the RHS generates
 2. how efficient the RHS access is (compare BW with max. BW)
 3. how much optimization potential we have with matrix reordering
- Do not forget about **load balancing!**
- Sparse matrix times **multiple vectors** bears the potential of huge savings in data volume
- **Consequence: Modeling is not always 100% predictive. It’s all about *learning more* about performance properties!**