

Performance Engineering – Case study:

Basic understanding of performance of sparse matrix vector multiplication (spMVM)

Prof. Dr. G. Wellein^(a,b) , Dr. G. Hager^(a), J. Hammer^(b), C.L. Alappat^(b)

^(a)HPC Services – Regionales Rechenzentrum Erlangen

^(b)Department für Informatik

University Erlangen-Nürnberg, Sommersemester 2019

Motivation – spMVM: most expensive part in many iterative solver – Conjugate Gradient Method



- Problem: Solve $\mathbf{A} \mathbf{x} = \mathbf{b}$ with \mathbf{A} large sparse ($n \times n$) matrix and \mathbf{b} is known vector
- Use (iterative) Conjugate Gradient (CG) Method

Example code in **MATLAB / GNU Octave** [edit]

```
function [x] = conjgrad(A, b, x)
    r = b - A * x;
    p = r;
    rsold = r' * r;

    for i = 1:length(b)
        Ap = A * p;
        alpha = rsold / (p' * Ap);
        x = x + alpha * p;
        r = r - alpha * Ap;
        rsnew = r' * r;
        if sqrt(rsnew) < 1e-10
            break;
        end
        p = r + (rsnew / rsold) * p;
        rsold = rsnew;
    end
end
```

Input vector b → `b`

vector r → `r = b - A * x;`

vector p → `p = r;`

result vector x → `[x]`

Sparse Matrix A → `A`

CG iteration → `for i = 1:length(b)`

vector Ap → `Ap = A * p;`

spMVM: Multiply A with vector p → `Ap = A * p;`

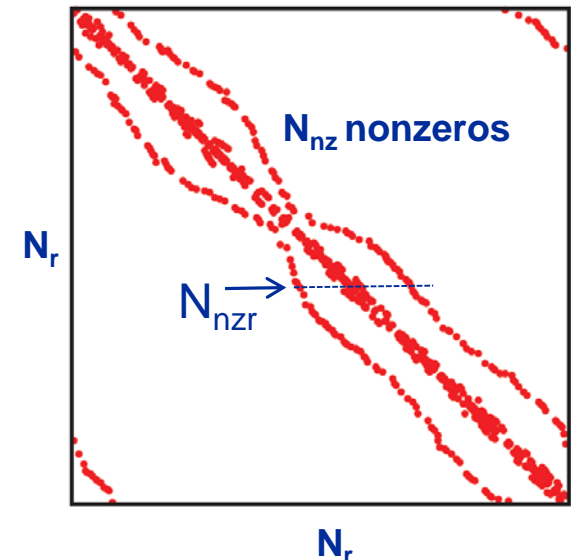
Vector-vector operations → `alpha = rsold / (p' * Ap);`
`r = r - alpha * Ap;`
`rsnew = r' * r;`

Stop iteration if converged → `if sqrt(rsnew) < 1e-10`
`break;`

https://en.wikipedia.org/wiki/Conjugate_gradient_method



- **Sparse Matrix Vector Multiply (spMVM) is the basic operation in many numerical applications, e.g.**
 - Finite-Element-Methods → Solve sparse linear system of equations
 - Quantum physics/chemistry → Determine eigenvalues of sparse matrices
- **Fraction of total time spent in spMVM is often approx. 85% – 99.99%**
- **“Sparse” matrix $\cong N_{nz}$ grows slower than quadratically with N_r**
 - N_{nz} = total number on nonzeros
 - N_{nzs} = avg. # nonzeros per row ($=N_{nz} / N_r$)
- **Store nonzero elements only:**
 - Reduce memory overhead
 - Avoid computing zeros



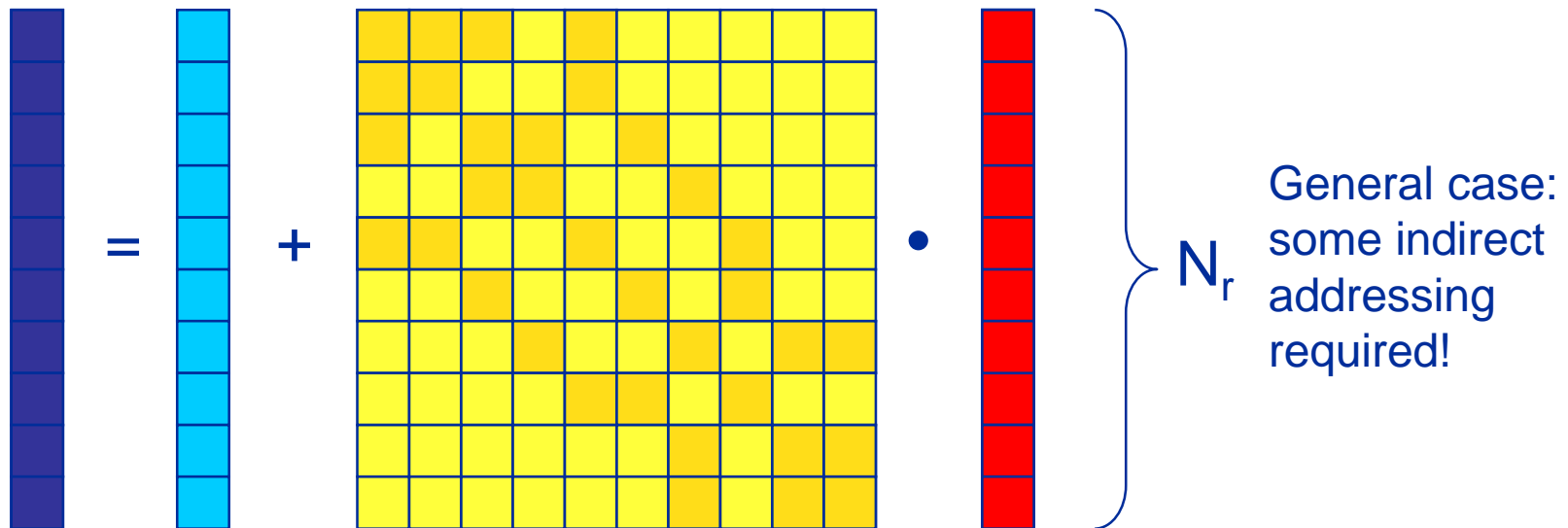
Sparse Matrix Vector Multiplication (spMVM)



Intensity-type $O(N)/O(N) \rightarrow$ memory bound

Nevertheless, there is more than one loop here!

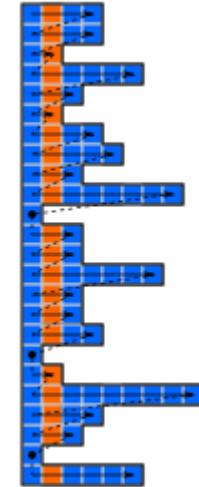
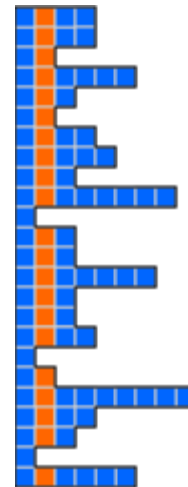
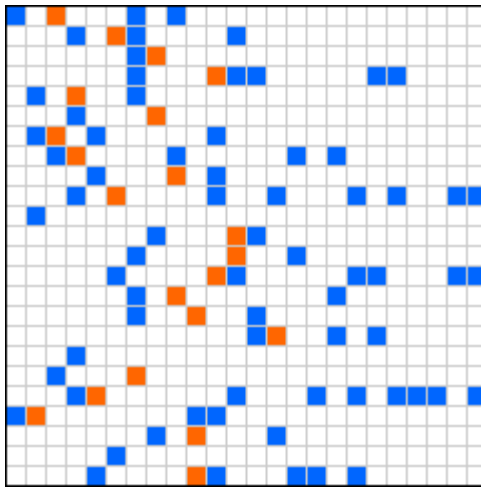
$$\mathbf{c} = \mathbf{c} + \text{sparse_mat} * \mathbf{b}$$



Sparsity pattern \leftrightarrow indirect access to right hand side vector (RHS \leftrightarrow vector **b** in our case)



- **Choice of sparse matrix storage scheme is crucial for performance**
 - Different schemes yield entirely different performance characteristics
 - Depends on architecture!
 - **Most important formats:**
 - **CRS (Compressed Row Storage)** → CPUs
 - (Sliced) ELLPACK → GPGPUs
 - **Other/newer possibilities:**
 - CCS (Compressed Column Storage, “Harwell-Boeing”)
 - CDS (Compressed Diagonal Storage)
 - **SELL-C-sigma** (Ours – perfect fit for SIMD & GPGPUs)
- M. Kreutzer et al: *A unified sparse matrix data format for efficient general sparse matrix-vector multiplication on modern processors with wide SIMD units*. SIAM SISC **36**(5), C401–C423 (2014). DOI: 10.1137/130930352
- **Depending on the storage scheme, the memory access patterns differ vastly between the formats**
 - Choose the storage scheme that best fits your needs

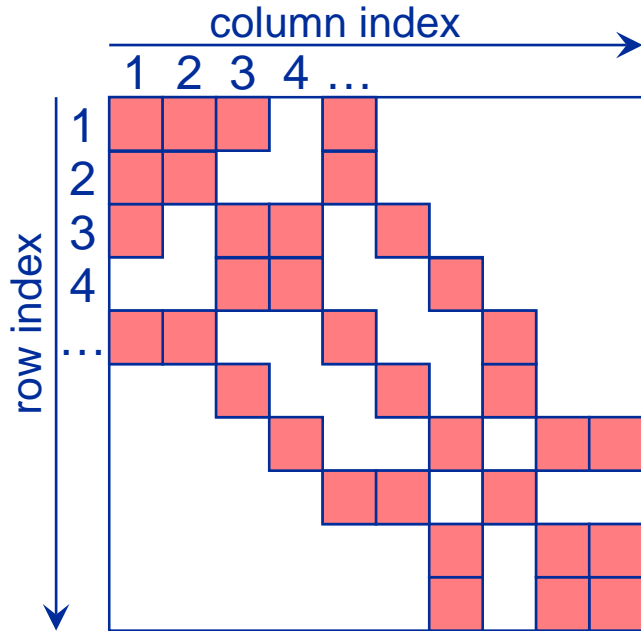


Format creation

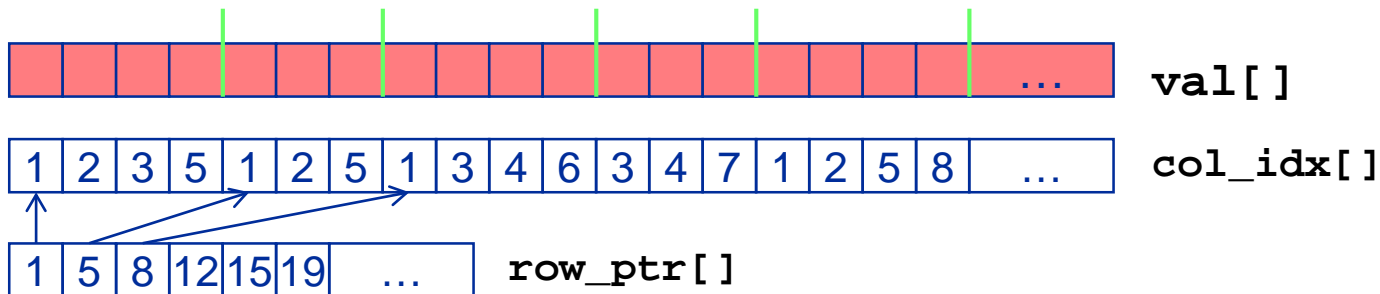
1. Store values and column indices of all non-zero elements **row-wise**
2. Store starting indices of each column (rpt)

Data fields

```
double val[]  
unsigned int col[]  
unsigned int rpt[]
```



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





$$A = \begin{pmatrix} 10 & 0 & 0 & 12 & 0 \\ 0 & 0 & 11 & 0 & 13 \\ 0 & 16 & 0 & 0 & 0 \\ 0 & 0 & 11 & 0 & 13 \end{pmatrix}$$

$$\begin{aligned} val &= \begin{pmatrix} 10 & 12 & 11 & 13 & 16 & 11 & 13 \\ (0,0) & (0,3) & (1,2) & (1,4) & (2,1) & (3,2) & (3,4) \end{pmatrix} \\ colInd &= \begin{pmatrix} 0 & 3 & 2 & 4 & 1 & 2 & 4 \\ (0) &) & (1) &) & (2) & (3) &) \end{pmatrix} \\ rowPtr &= \begin{pmatrix} 0 & 2 & 4 & 5 & 7 \\ (0) & (1) & (2) & (3) & (4) \end{pmatrix} \end{aligned}$$

- Only the necessary (= non-zero) values along with their column indices are stored



- Implement $\mathbf{c} = \mathbf{c} + \text{sparse_mat} * \mathbf{b}$
- Only the nonzero elements of the matrix are stored/used
 - Operation count = $2 * N_{nz}$ [FLOP]
 - Data set size for matrix storage: $(8 + 4) * N_{nz}$ [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
```

- Features
 - **Long outer loop** (N_r) – matrix dimension: $N_r \sim 10^4, \dots, 10^{10}$
 - Probably **short inner loop** ($\sim N_{nzs}$): $N_{nzs} = N_{nz} / N_r \sim 10s, \dots, 100s$
(nonzero entries in each respective row)
 - **Indexed (indirect) access** to RHS vector $\mathbf{b}[\]$
 - All other data is accessed contiguously!



```
!$OMP PARALLEL DO SCHEDULE(runtime)
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
```

- OMP_SCHEDULE = STATIC
 - OK for “balanced” matrix structures
 - Otherwise potential load imbalance
- Other schedules may achieve better performance

Performance characteristics

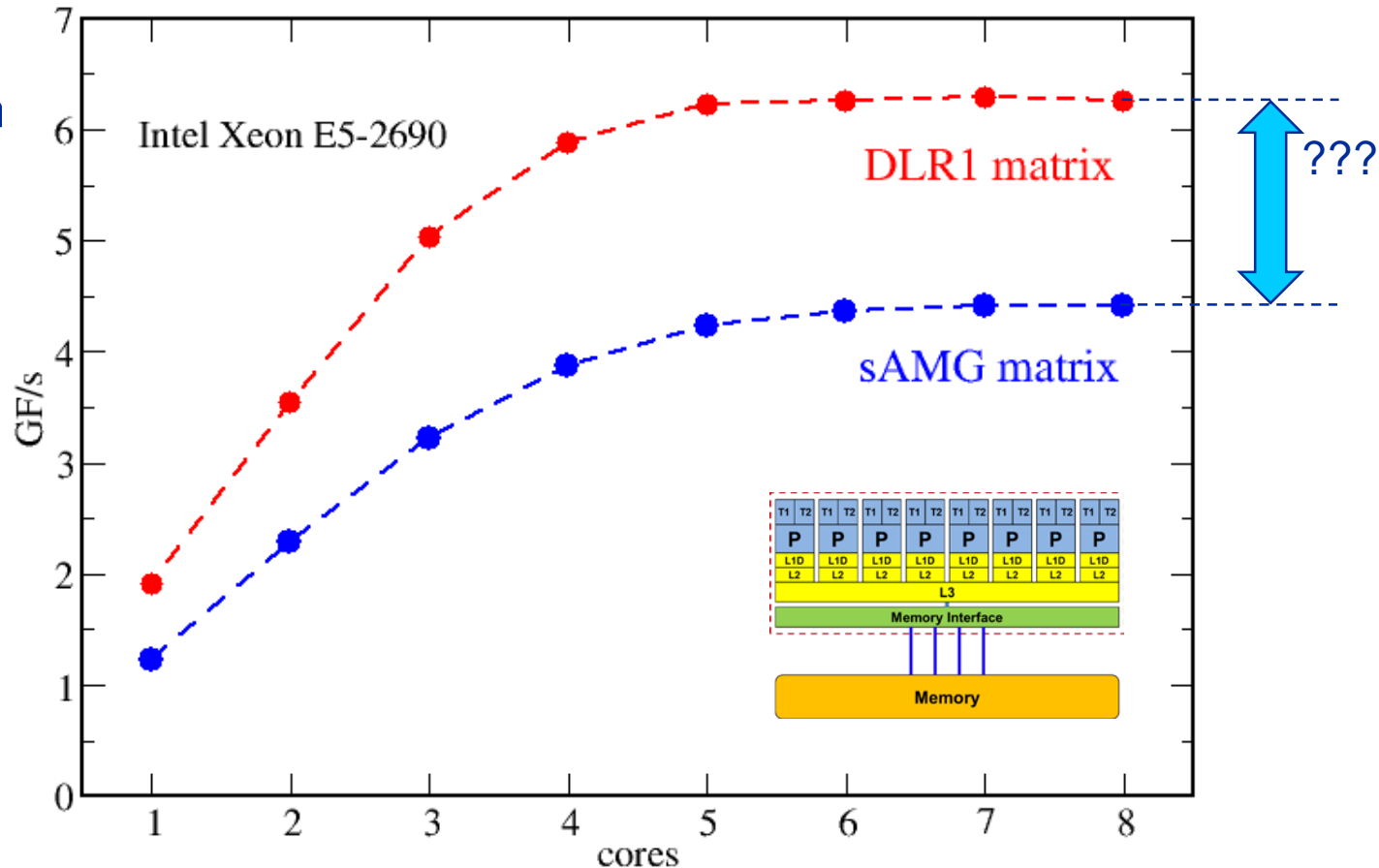


- Strongly memory-bound for large data sets → saturating performance across cores on the chip
- Performance seems to depend on the matrix

Can we explain this?

Is there a “light speed” for spMVM?

Optimization?





- **Sparse MVM in double precision w/ CRS data storage:**

```

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
    
```

- **DP CRS comp. intensity**

$$I_{CRS}^{DP} = \frac{2}{8 + 4 + 8\alpha + 16/N_{nzs}} \frac{\text{flops}}{\text{byte}}$$

- α quantifies traffic for loading RHS (depends on matrix!!!!)
 - $\alpha = 0 \rightarrow$ RHS is in cache
 - $\alpha = 1/N_{nzs} \rightarrow$ RHS loaded once
 - $\alpha = 1 \rightarrow$ no cache
 - $\alpha > 1 \rightarrow$ Houston, we have a problem!
- “Expected” performance = $b_s \times I_{CRS}$
- Determine α by measuring performance and actual memory traffic
 - Maximum memory BW may not be achieved with spMVM



$$I_{CRS}^{DP} = \frac{2}{8 + 4 + 8\alpha + 16/N_{nzs}} \frac{\text{flops}}{\text{byte}} = \frac{N_{nz} \cdot 2 \text{ flops}}{V_{meas}}$$

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

- **Solve for α :**

$$\alpha = \frac{1}{4} \left(\frac{V_{meas}}{N_{nz} \cdot 2 \text{ bytes}} - 6 - \frac{8}{N_{nzs}} \right)$$

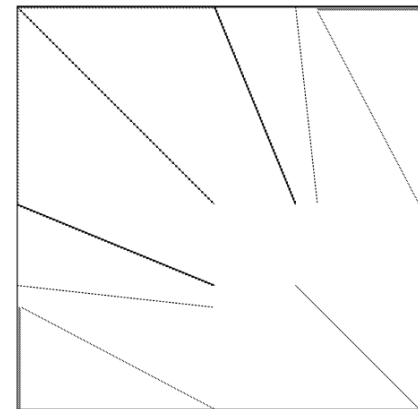
- **Example: kkt_power matrix from the UoF collection on one Intel SNB socket**

- $N_{nz} = 14.6 \cdot 10^6, N_{nzs} = 7.1$
- $V_{meas} \approx 258 \text{ MB}$
- $\rightarrow \alpha = 0.43, \alpha N_{nzs} = 3.1$
- \rightarrow **RHS is loaded 3.1 times** from memory

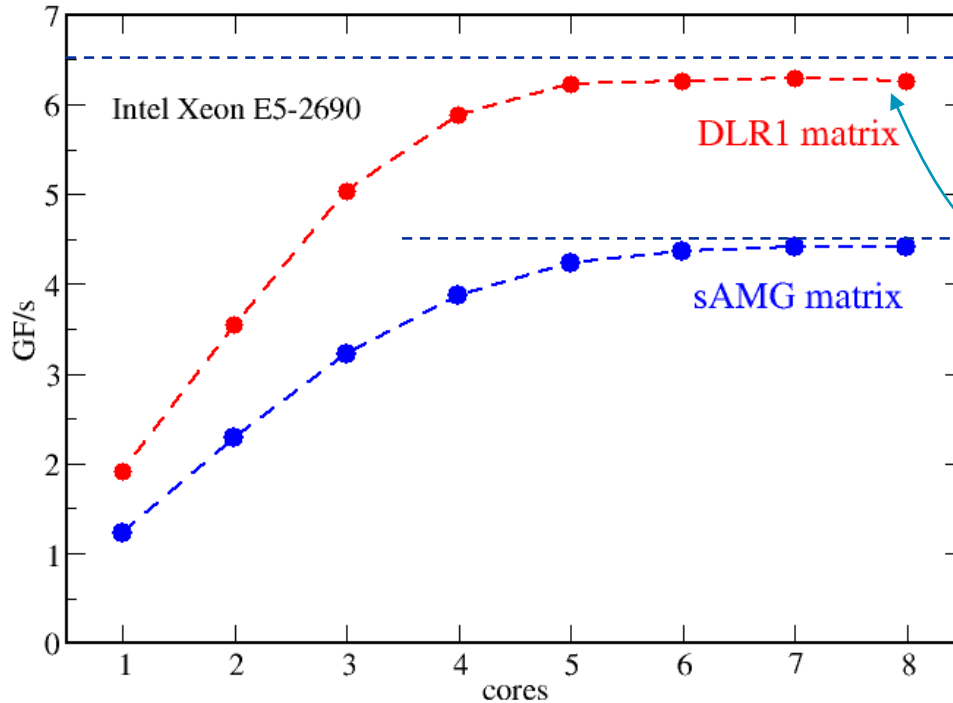
- and:

$$\frac{I_{CRS}^{DP}(1/N_{nzs})}{I_{CRS}^{DP}(\alpha)} = 1.15$$

15% extra traffic \rightarrow optimization potential!



Now back to the start...



- $b_S = 39 \text{ GB/s}$
- $B_c^{min} = 6 \text{ B/F}$
- Maximum spMVM performance:

$$P_{max} = 6.5 \text{ GF/s}$$

- → DLR1 causes minimum code balance!
- sAMG matrix code balance:

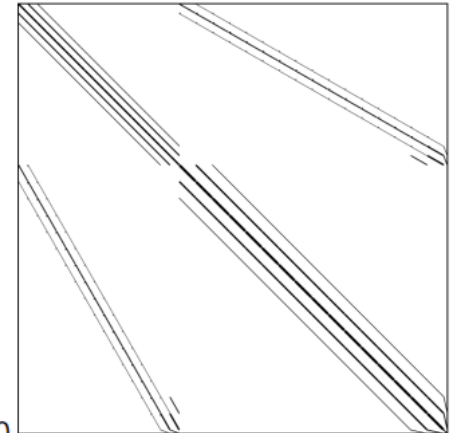
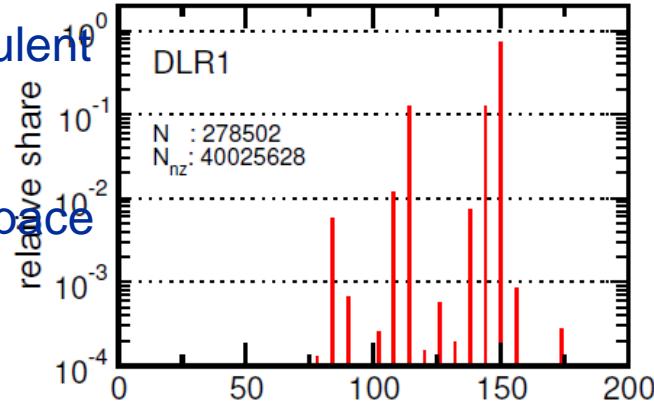
$$B_c \leq \frac{b_S}{4.5 \text{ GF/s}} = 8.7 \text{ B/F}$$

- Why is this only an upper limit?
- What is the next step?
- Could we have predicted this qualitative difference?



“DLR1” (A. Basermann, DLR)

Adjoint problem computation (turbulent transonic flow over a wing) with the TAU CFD system of the German Aerospace Center (DLR)
 Avg. non-zeros/row ~150

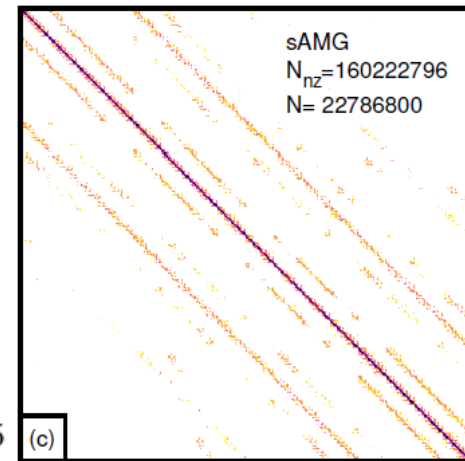
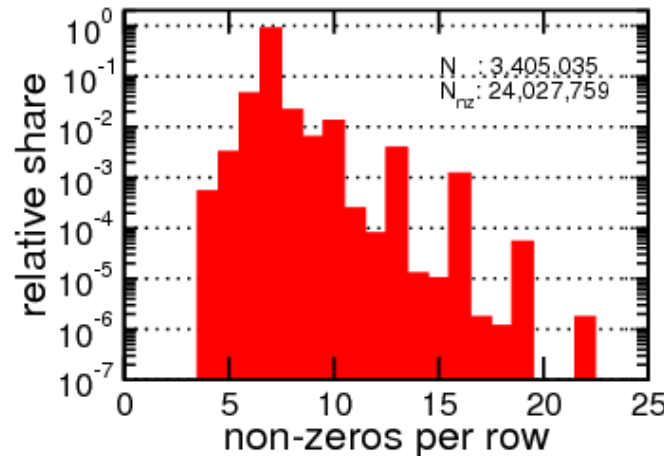


“sAMG” (K. Stüben, FhG-SCAI)

Matrix from FhG’s adaptive multigrid code sAMG for the irregular discretization of a Poisson problem on a car geometry.
 Avg. non-zeros/row: $N_{nzs} \sim 7$

Assume $\alpha = \frac{1}{N_{nzs}} = \frac{1}{7}$

$$I_{CRS}^{DP} = \frac{1}{2} \frac{\text{flops}}{12 + 8/7 + 16/7 \text{ byte}} = \frac{1}{8} \frac{F}{B}$$





- **Conclusion from Roofline analysis**
 - The roofline model does not work 100% 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
- **Consequence: Modeling is not always 100% predictive. It's all about *learning more* about performance properties!**