



Erlangen Regional
Computing Center



Case study: Sparse Matrix-Vector Multiplication



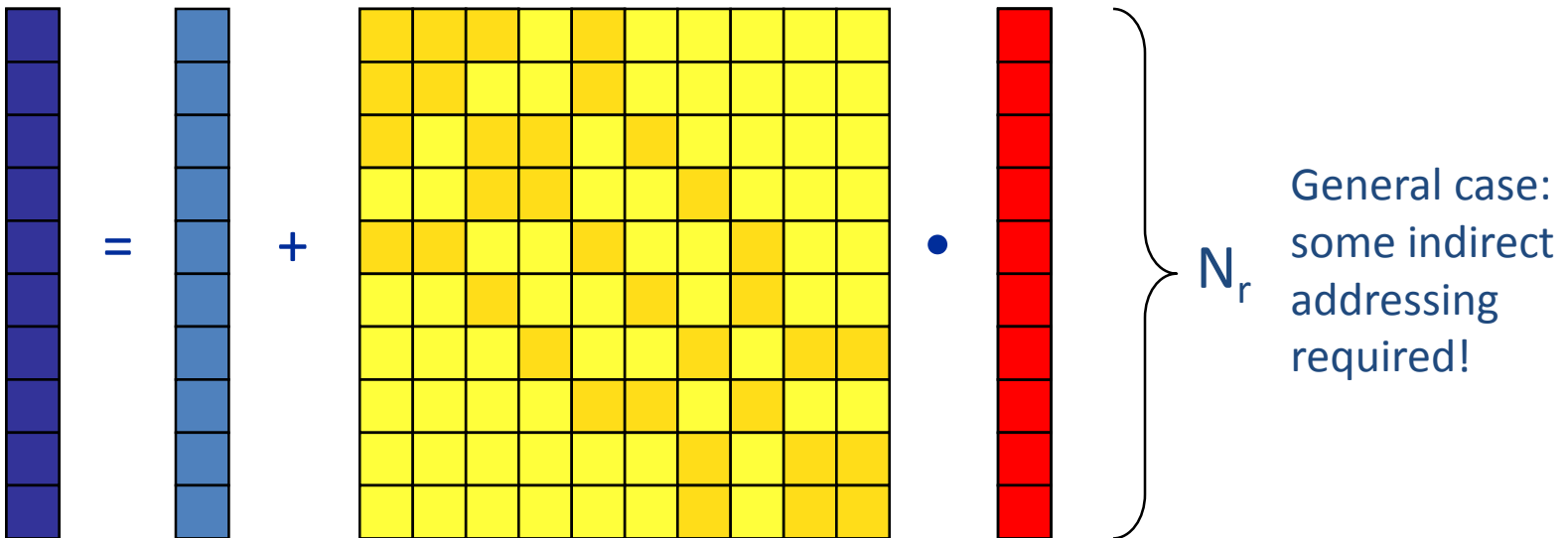


Erlangen Regional
Computing Center

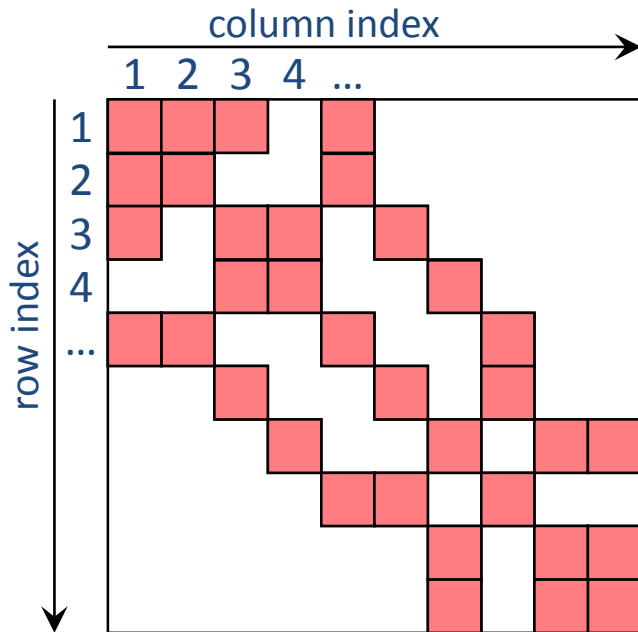


SpMVM: The Basics

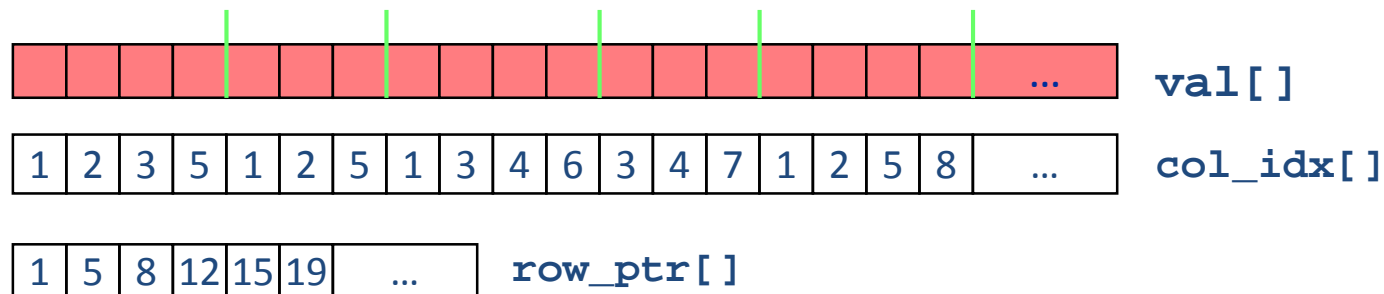
- Key ingredient in some matrix diagonalization algorithms
 - Lanczos, Davidson, Jacobi-Davidson
- Store only N_{nz} nonzero elements of matrix and RHS, LHS vectors with N_r (number of matrix rows) entries
- “Sparse”: $N_{nz} \sim N_r$
- Average number of nonzeros per row: $N_{nzs} = N_{nz}/N_r$



- 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...



Erlangen Regional
Computing Center



SpMVM: Performance Analysis

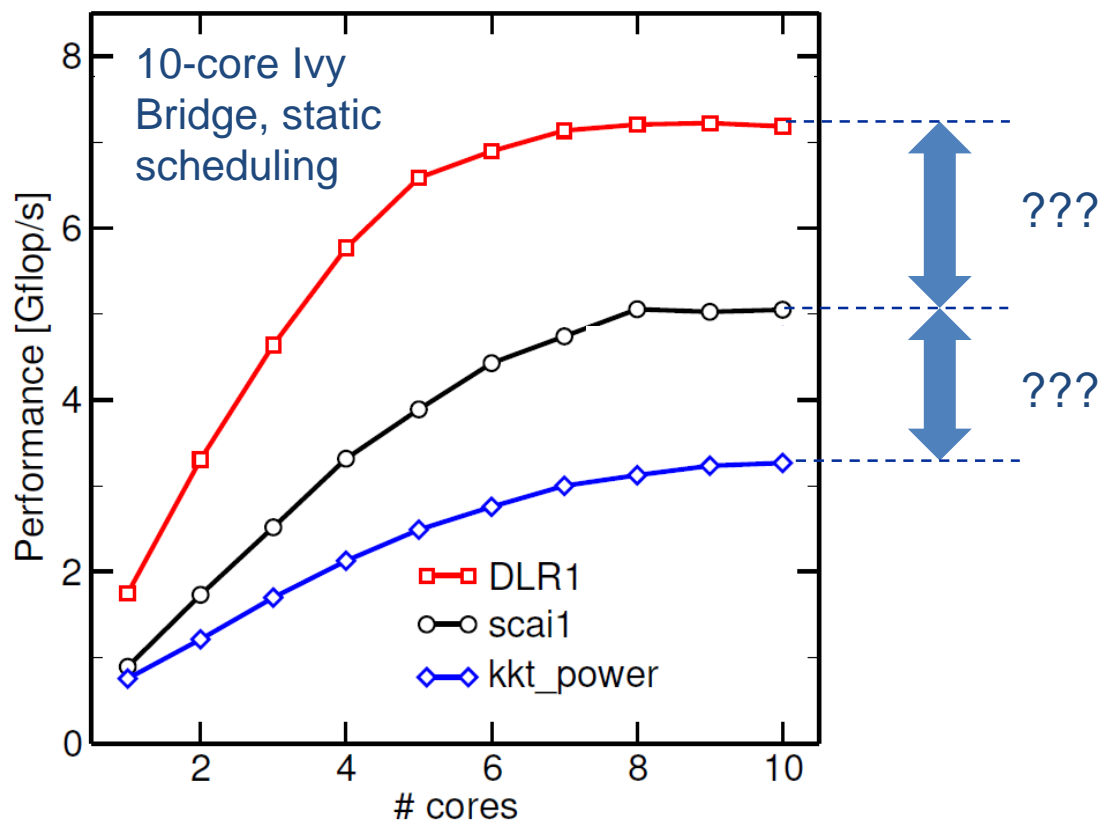


- Strongly memory-bound for large data sets → saturating performance across cores on the chip
- Performance seems to depend on the matrix

- Can we explain this?

- Is there a “light speed” for SpMV?

- 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
```

$$B_c^{DP, CRS} = \frac{8 + 4 + 8\alpha + 20/N_{nzs}}{2} \frac{B}{F} = \left(6 + 4\alpha + \frac{10}{N_{nzs}} \right) \frac{B}{F}$$

Absolute minimum code balance: $B_{\min} = 6 \frac{B}{F}$
 $\rightarrow I_{\max} = \frac{1}{6} \frac{F}{B}$



Hard upper limit for
in-memory
performance: b_s/B_{\min}

DP CRS code balance

- α quantifies the traffic for loading the RHS
 - $\alpha = 0 \rightarrow$ 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!
- “Target” performance = b_S/B_c
- **Caveat:** Maximum memory BW may not be achieved with spMVM (see later)

$$B_c^{DP,CRS}(\alpha) = \frac{8+4+8\alpha+20/N_{nzs}}{2} \frac{B}{F}$$
$$= \left(6+4\alpha+\frac{10}{N_{nzs}}\right) \frac{B}{F}$$

Can we predict α ?

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

\rightarrow Determine α by measuring **the actual memory traffic**

Determine α (RHS traffic quantification)

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

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

- Solve for α :

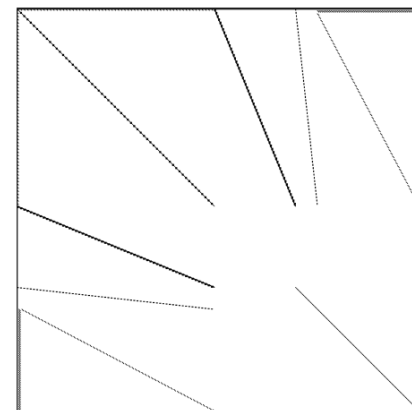
$$\alpha = \frac{1}{4} \left(\frac{V_{meas}}{N_{nzs} \cdot 2 \text{ bytes}} - 6 - \frac{10}{N_{nzs}} \right)$$

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:

$$\frac{B_c^{DP,CRS}(\alpha)}{B_c^{DP,CRS}(1/N_{nzs})} = 1.11$$

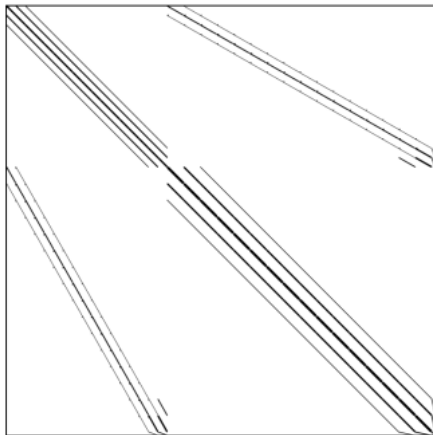
11% extra traffic \rightarrow
optimization potential!



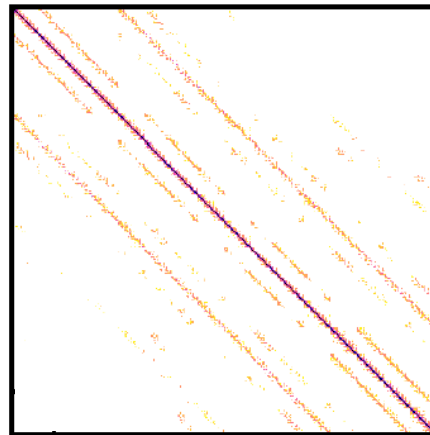
Three different sparse matrices

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

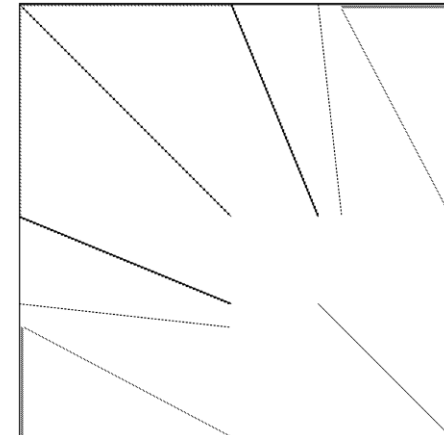
Matrix	N	N_{nzs}	B_c^{opt} [B/F]	P_{opt} [GF/s]
DLR1	278,502	143	6.1	7.64
scai1	3,405,035	7.0	8.0	5.83
kkt_power	2,063,494	7.08	8.0	5.83



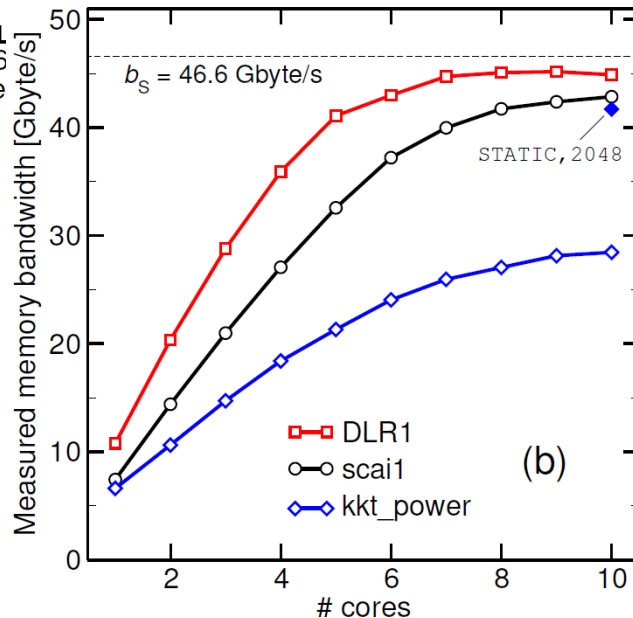
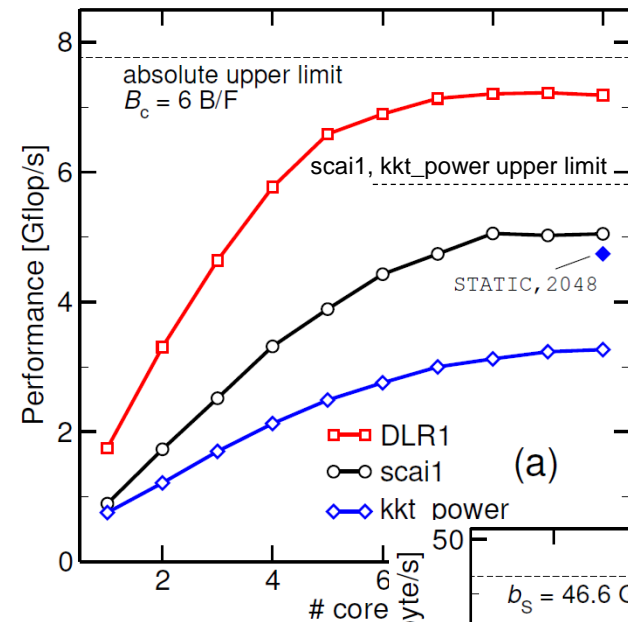
DLR1



scai1



kkt_power



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

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

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

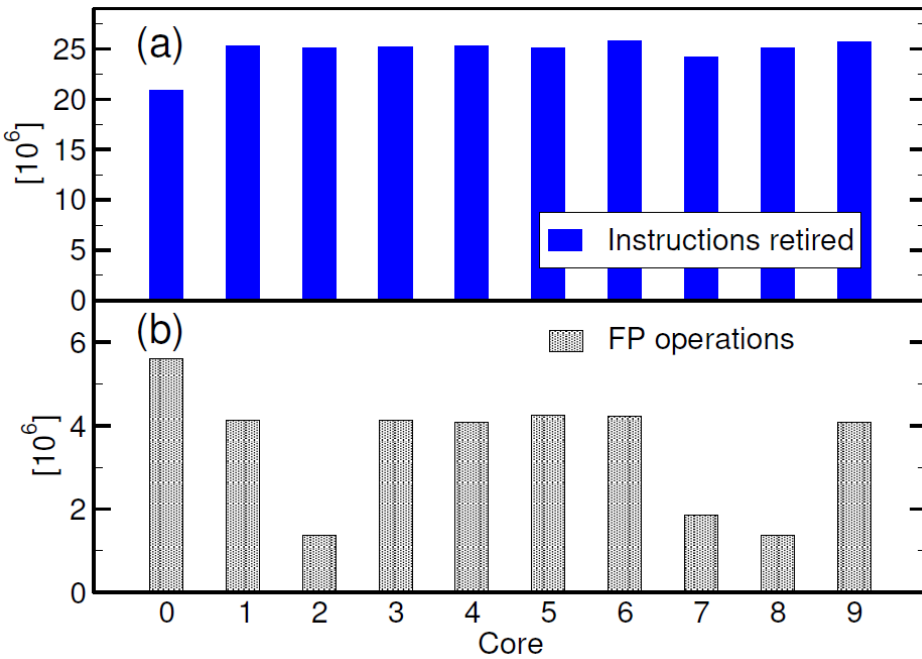
$$B_c^{meas} \approx 8.5 \text{ B/F} > B_c^{opt}$$

- → good BW utilization, slightly non-optimal α

- **kkt_power** measured balance:

$$B_c^{meas} \approx 8.8 \text{ B/F} > B_c^{opt}$$

- → performance degraded by load imbalance, fix by block-cyclic schedule

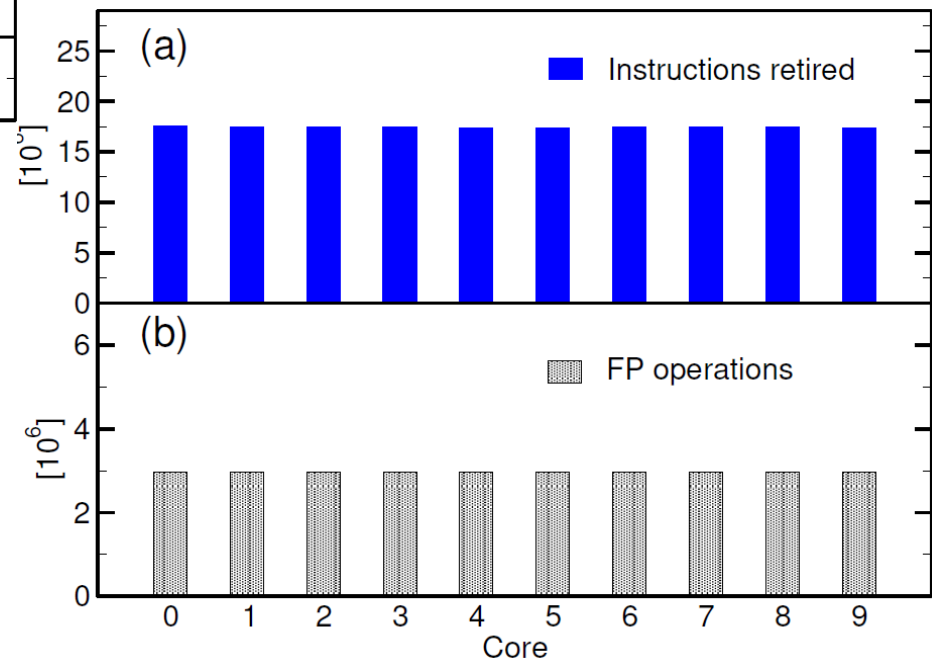


Measurements with likwid-perfctr (MEM_DP group)



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

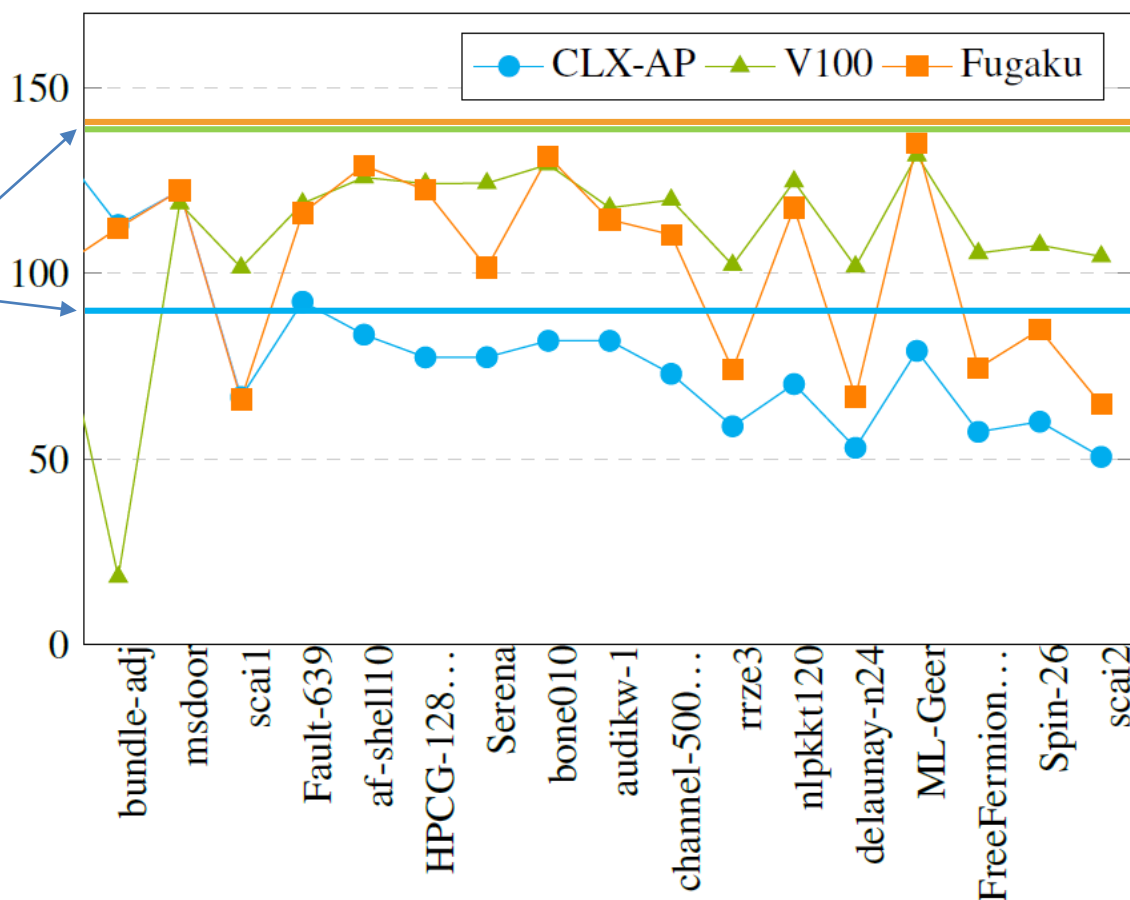
→ CPI value unchanged!



CPU performance comparison

- Cascade-Lake AP: $b_S = 500$ Gbyte/s
- NVIDIA V100: $b_S = 840$ Gbyte/s
- Fujitsu A64FX in Fugaku: $b_S = 859$ Gbyte/s

Absolute upper limits
(matrix independent)
given by $b_S / (6 \text{ B/F})$





Erlangen Regional
Computing Center



SpMVM with multiple RHS & LHS Vectors



Multiple RHS vectors (SpMMV)

Unchanged matrix applied to multiple RHS vectors to yield multiple LHS 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
```

Lower 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!

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

$N_{nzs} \gg 1$

$$\frac{6}{r} \frac{B}{F}$$

$r = 1$

$$\left(6 + 4\alpha + \frac{10}{N_{nzs}} \right) \frac{B}{F}$$

reassuring 😊

$r \gg 1$

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

Can become very small for large $N_{nzs} \rightarrow$ decoupling from memory bandwidth is possible!

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!**