

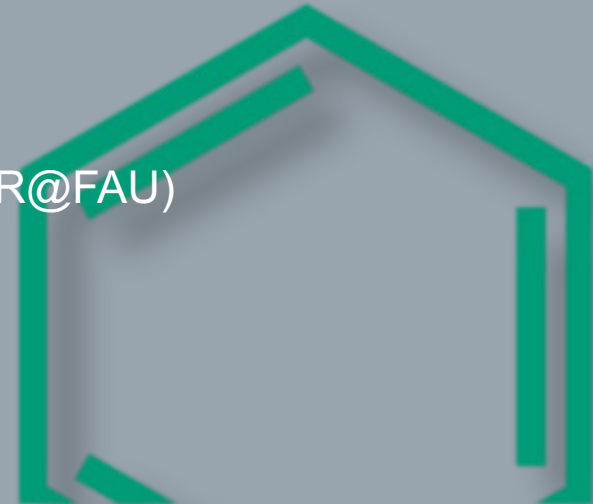
Performance Engineering for Sparse Matrix-Vector Multiplication: Some new ideas for old problems

Christie Alappat, Georg Hager, Gerhard Wellein

Erlangen National High Performance Computing Center (NHR@FAU)

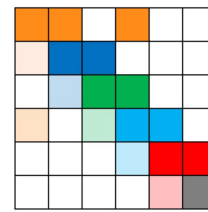
Department of Computer Science

Friedrich-Alexander-Universität Erlangen-Nürnberg



Agenda

- Motivation
- RACE: A different approach to Sparse Matrix-Vector Multiplication (SpMV)
- Parallelization of Symmetric SpMV
- Cache blocking for Matrix Power Kernels
- Conclusion

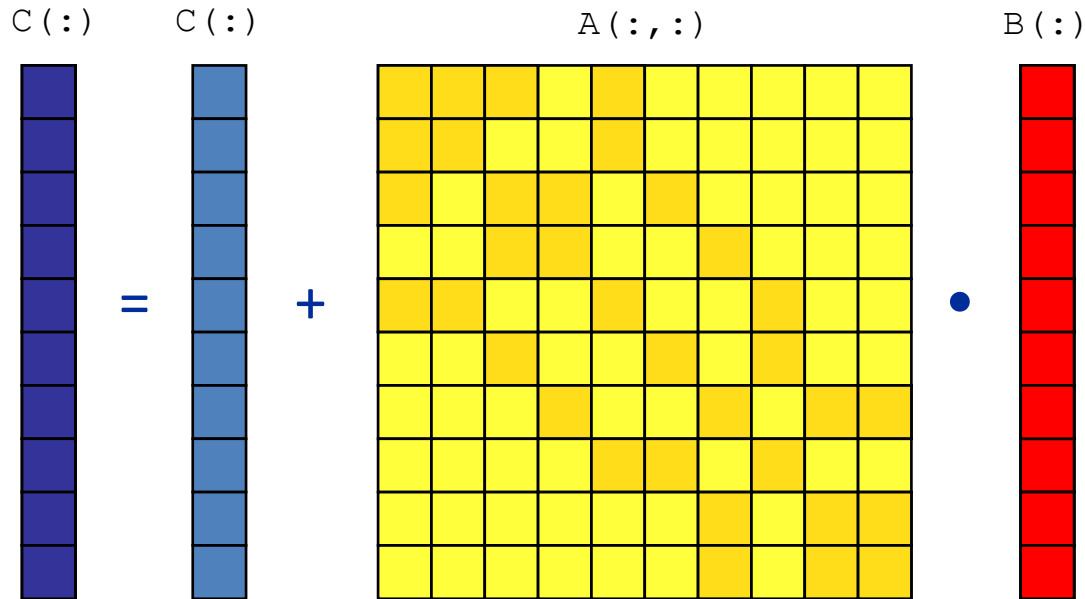


$$y = A^p x$$

!WARNING! This talk considers single multicore processors only:
Cascade Lake, Ice Lake, Rome (20c – 64c)

Motivation – Sparse Matrix Vector Multiplication

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



Performance of SpMV – Basics: Storage formats

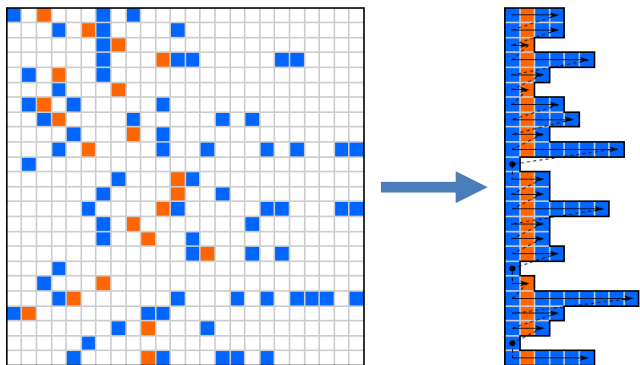
- Store and use non-zero entries only
- Hardware-efficient formats for storing non-zeros in $\mathbb{A}(:, :)$
 - **CRS**, COO, DIA,...
 - ELLPACK, hybrid,... → GPU formats, e.g. [1]
 - SELL-C- σ , SELL-P,...
- Best choice \leftrightarrow Matrix structure & Architecture

Goals:

- Consecutive access to $\mathbb{A}(:, :)$
- Low storage overhea

[1] S. Filippone et al.: Sparse Matrix-Vector Multiplication on GPGPUs. ACM TOMS (2017) <https://doi.org/10.1145/3017994>

Performance of SpMV – Basics: Performance Model



```

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

} A(:, :)

```

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
    
```

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

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

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

Motivation – Sparse Matrix Vector Multiplication

- Improve code balance of **SpMV-algorithms**

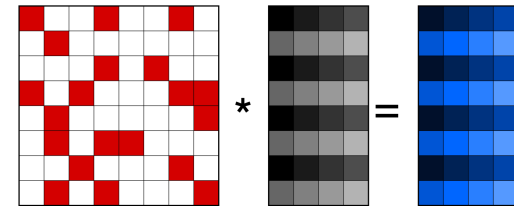
- Kernel fusion, e.g. HPCG:

Algorithm 2. HPCG

```
1: while  $k \leq \text{iter}$  &  $r_{\text{norm}}/r_0 > \text{tol}$  do
2:    $z = MG(A,r)$            -> MG sweep
3:    $\text{oldrtz} = \text{rtz}$ 
4:    $\text{rtz} = \langle r,z \rangle$        -> DOT
5:    $\beta = \text{rtz}/\text{oldrtz}$ 
6:    $p = \beta * p + z$        -> WAXPY
7:    $Ap = A * p$            -> SpMPV
8:    $pAp = \langle p, Ap \rangle$    -> DOT
9:    $\alpha = \text{rtz}/pAp$ 
10:   $x = x + \alpha * p$        -> WAXPY
11:   $r = r - \alpha * Ap$      -> WAXPY
12:   $r_{\text{norm}} = \langle r,r \rangle$    -> DOT
13:   $r_{\text{norm}} = \text{sqrt}(r_{\text{norm}})$ 
14:   $k++$ 
```

- Sparse Matrix Multiple Vector Multiplication

Gropp et al.: Towards Realistic Performance Bounds for Implicit CFD Codes
ParCFD 1999. <https://wgropp.cs.illinois.edu/bib/papers/pdata/1999/pcf99/gkks.ps>



- Both, e.g. Chebyshev Filter Computation for eigenvalue computations:

Kreutzer et al.: Chebyshev Filter Diagonalization on Modern Manycore Processors and GPGPUs. ISC 2018.
DOI: https://dx.doi.org/10.1007/978-3-319-92040-5_17

Motivation – Sparse Matrix Vector Multiplication

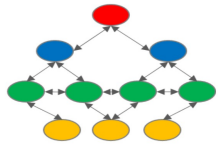
Further opportunities to improve code balance of SpMV-kernels:

- Exploit symmetry (SymmSpMV): $A = A^t$
 - Efficient parallelisation \leftrightarrow Write conflicts
 - Naive speed-up: max. 2

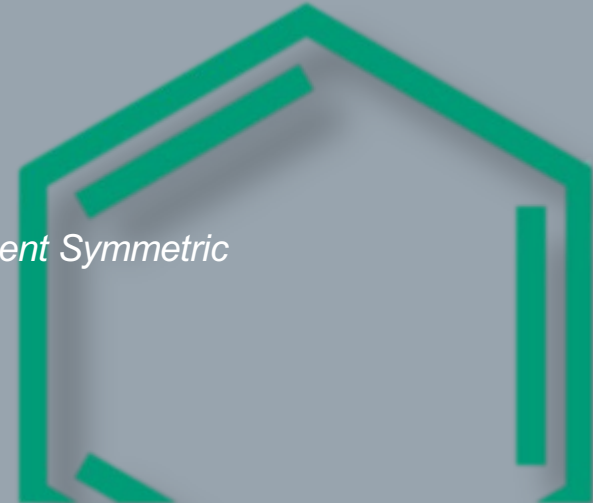
- Cache blocking for Sparse **Matrix Power Kernels (MPK)**: $y = A^p x$
 - Data locality in „sparse matrix-matrix multiply“ \leftrightarrow Irregular sparsity pattern
 - Naive speed up: max. p

SpMV – graph traversal – RACE

RACE
Hardware Friendly Coloring



Alappat et al., *A Recursive Algebraic Coloring Technique for Hardware-efficient Symmetric Sparse Matrix-vector Multiplication*. ACM Trans. Parallel Comput., 2020, DOI:10.1145/3399732



Consider SpMV as „graph traversal“

- Standard view: Run over all rows $i=1, \dots, n_r$ and calculate

$$y_i = \sum_{j \in \text{COL}(i)} A_{i,j} x_j$$

where $\text{COL}(i)$ contains the column indices of the non-zeros in i -th row

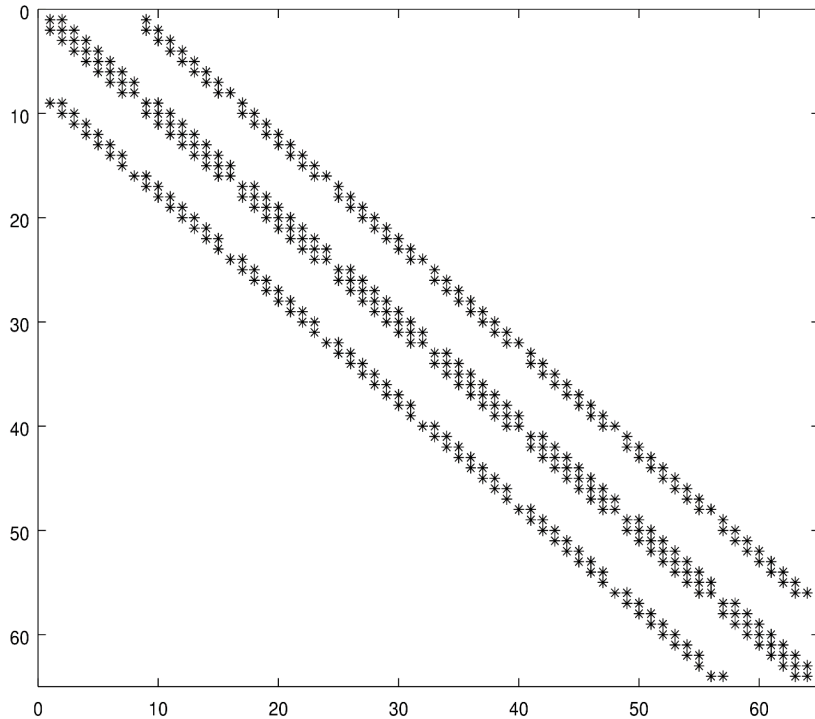
- RACE – **graph-based view**: row \leftrightarrow vertex & non-zero entry \leftrightarrow edge

In the graph terminology a SpMV operation ($y = Ax$) can be formulated as follows: If $G = (V, E)$ is the graph representation of the sparse matrix A then for every vertex $u \in V(G)$ calculate

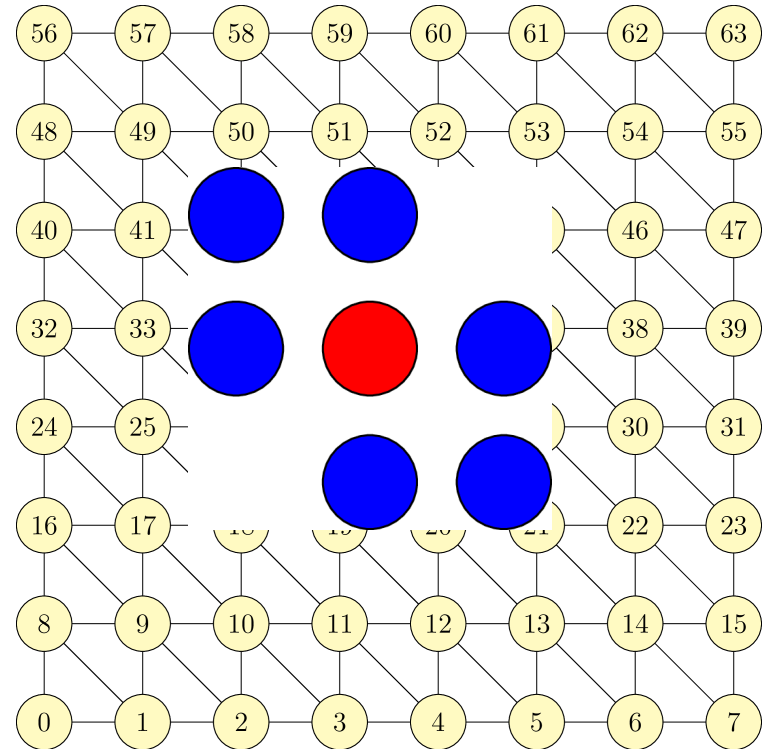
$$y_u = \sum_{v \in N(u)} A_{u,v} x_v . \quad (2)$$

Sample Stencil Matrix and its graph representation

Symmetric Matrix (Stencil)

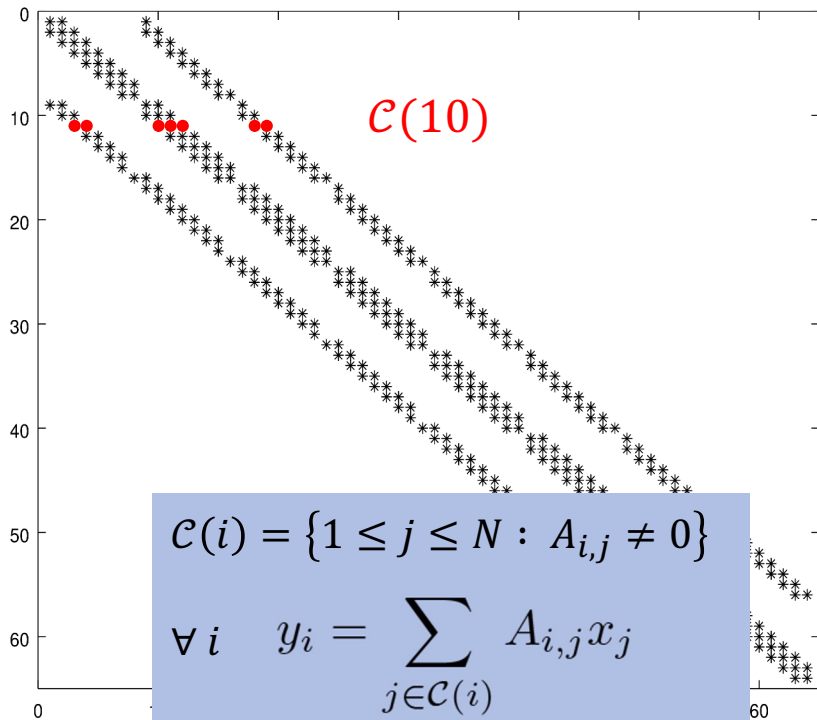


Undirected Graph

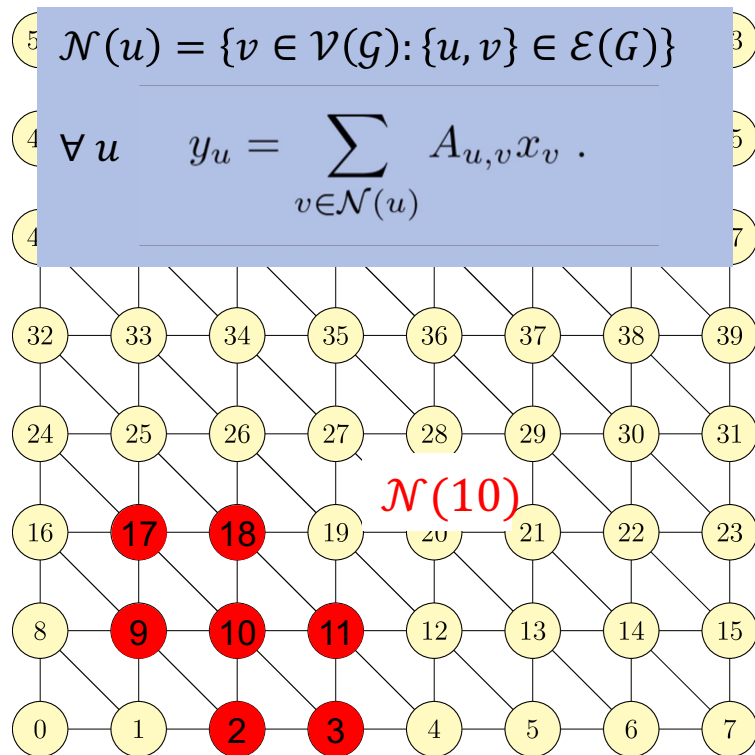


Sample Stencil Matrix and its graph representation

Symmetric Matrix (Stencil)

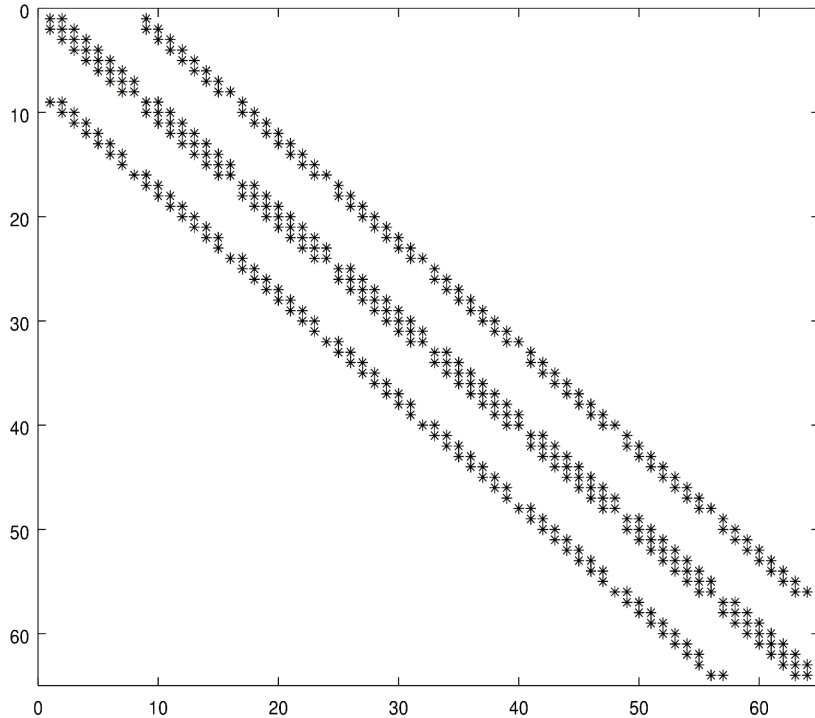


Undirected Graph

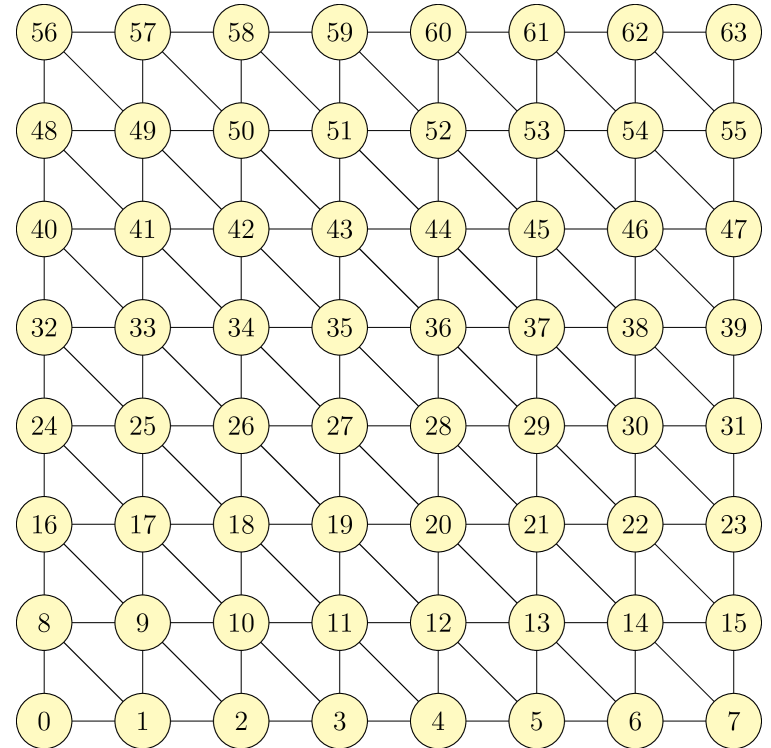


Sample Stencil Matrix and its graph representation

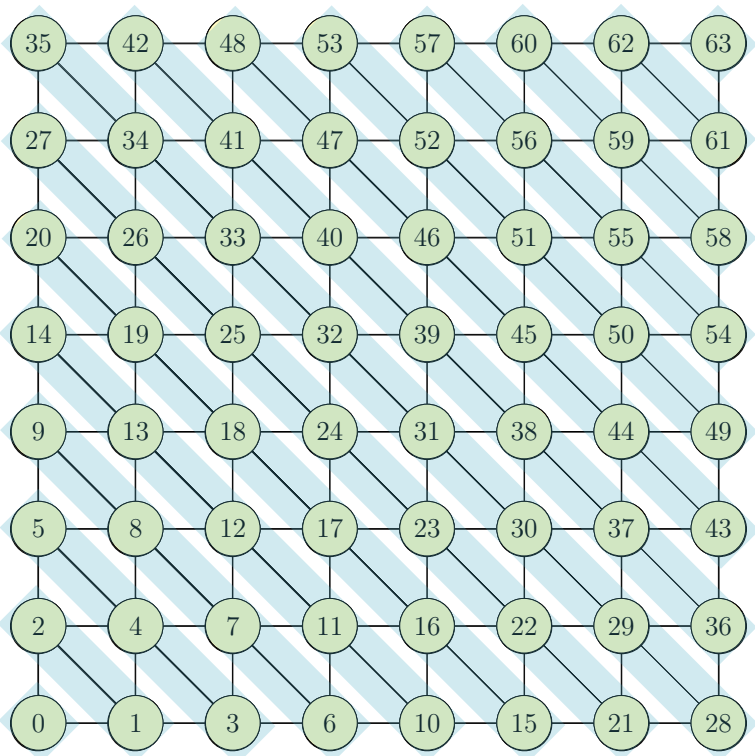
Symmetric Matrix (Stencil)



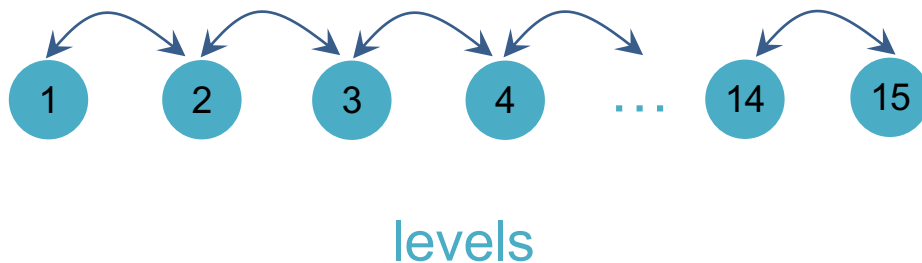
Undirected Graph



RACE: Get BFS-levels

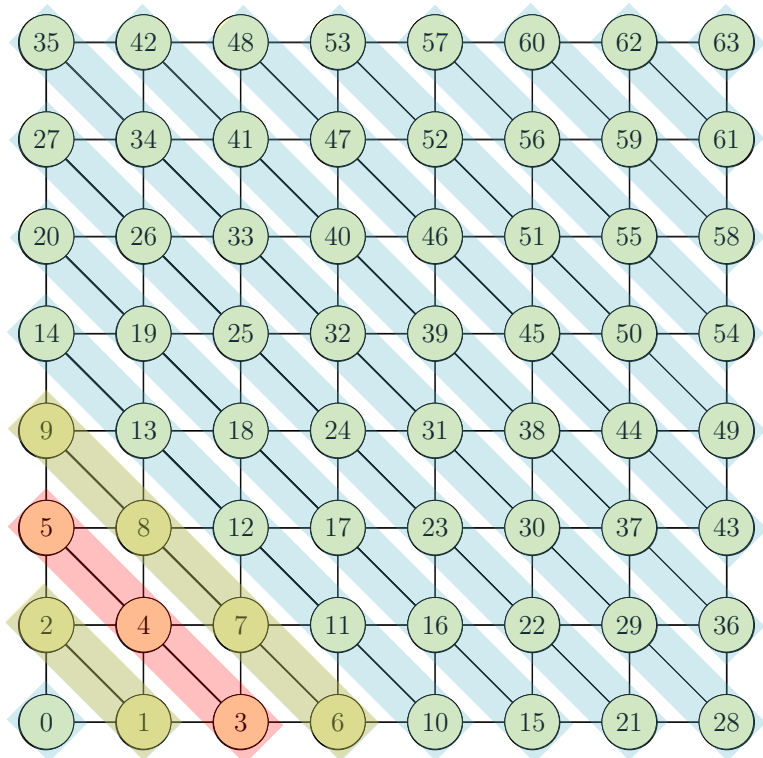


- Get levels of a Breadth-First-Search (BFS)
- Reorder vertices \rightarrow consecutive within level

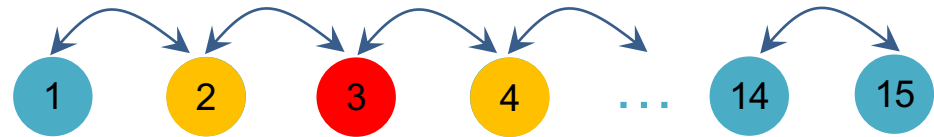


- Dependencies only between neighboring levels!

RACE: SpMV - Dependencies and Implementation



Update 3 → indirect access to 3, 2 and 4

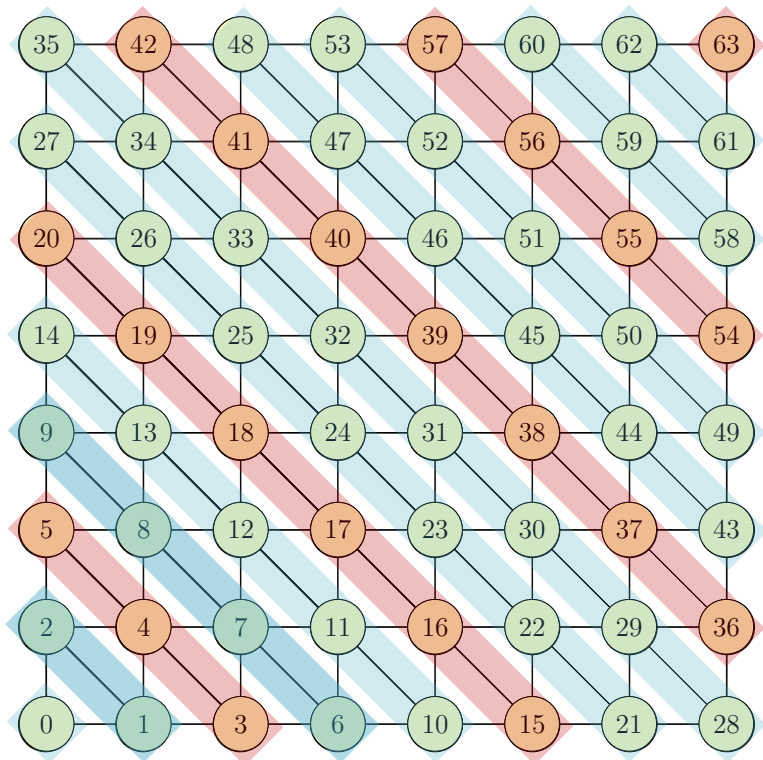


levels

```
do i = 1,L //loop over Levels
  SpMV_CRS(level_ptr[i], level_ptr[i+1])
enddo
```

```
function SpMV_CRS(start, end)
  do i = start, end
    do j = row_ptr(i), row_ptr(i+1) - 1
      y(i) = y(i) + val(j) * y(col_idx(j))
    enddo
  enddo
enddo
```

RACE: SymmSpMV - Parallelization



Symmetric SpMV – basic idea:
Compute distance-2 levels in parallel



C. Alappat, A. Basermann, A. R. Bishop, H. Fehske, G. Hager, O. Schenk, J. Thies, and G. Wellein:

A Recursive Algebraic Coloring Technique for Hardware-efficient Symmetric Sparse Matrix-vector Multiplication.

ACM Trans. Parallel Comput. (2020). DOI: 10.1145/3399732

RACE & Cache Blocking for MPK

IEEE TRANSACTIONS ON PARALLEL AND DISTRIBUTED SYSTEMS, VOL. 34, NO. 2, FEBRUARY 2023

581

Level-Based Blocking for Sparse Matrices: Sparse Matrix-Power-Vector Multiplication

Christie Alappat , Georg Hager , Olaf Schenk , *Senior Member, IEEE*, and Gerhard Wellein 

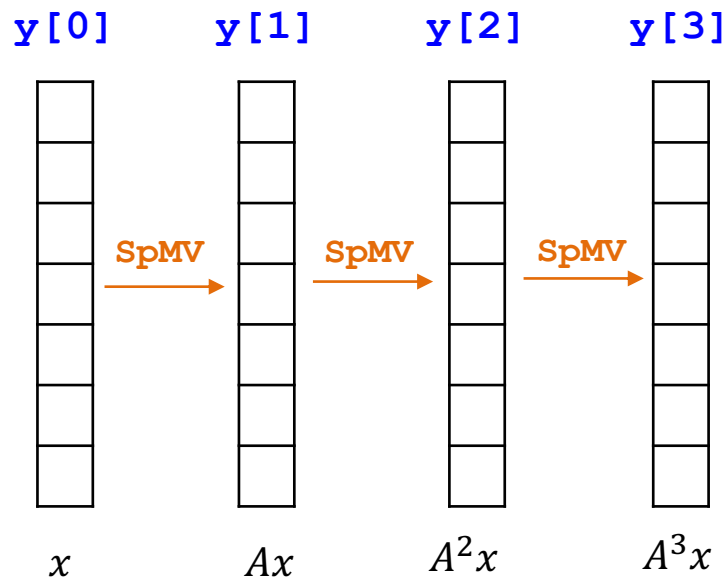
C. Alappat, et al, doi: 10.1109/TPDS.2022.3223512



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?

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!

RACE

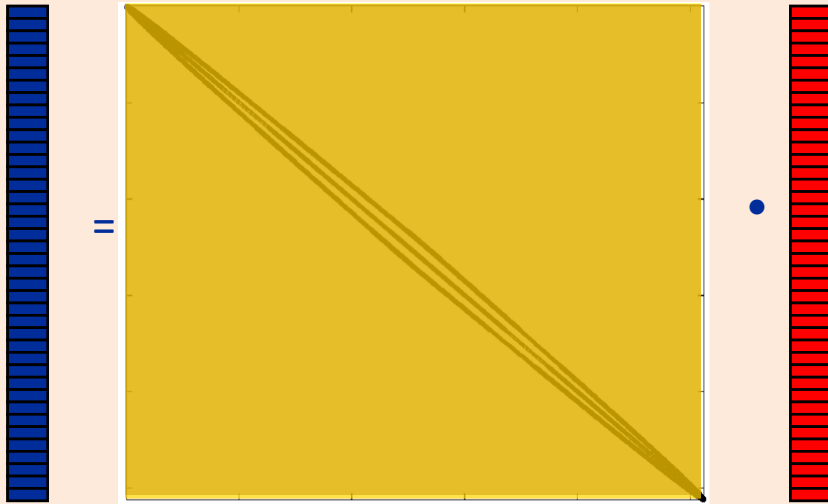
¹Huber et al., 2021. Graph-based multi-core higher-order time integration of linear autonomous partial differential equations. J. Comput. Sci. [DOI:10.1016/j.jocs.2021.101349](https://doi.org/10.1016/j.jocs.2021.101349)

²Mohiyuddin et al., 2009. Minimizing communication in sparse matrix solvers. In Proceedings of the SC'09. [DOI:10.1145/1654059.1654096](https://doi.org/10.1145/1654059.1654096)

Matrix power – Traditional approach vs. Cache Blocking

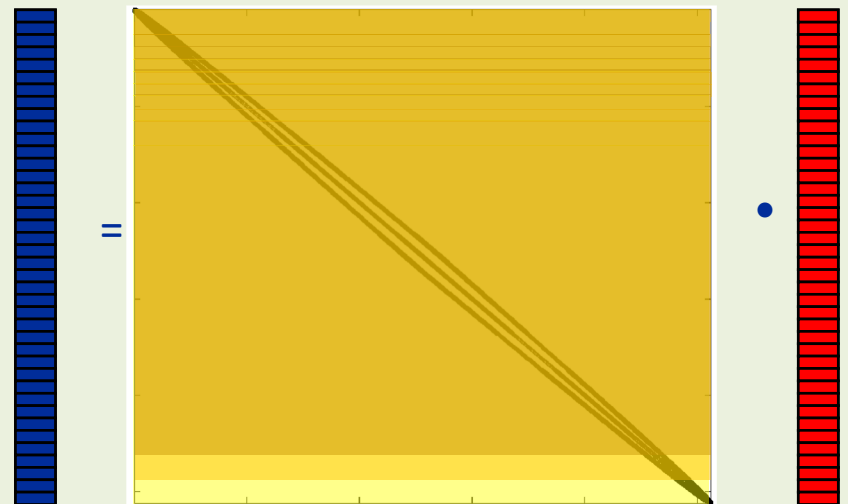
Calculate $y = A^3x$

TRAD approach



Matrix accessed 3 times from memory

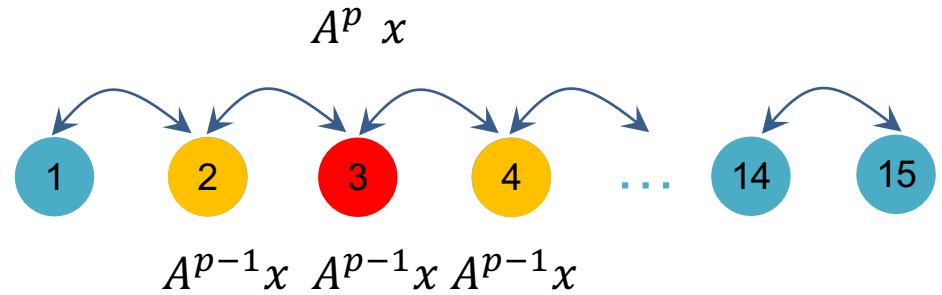
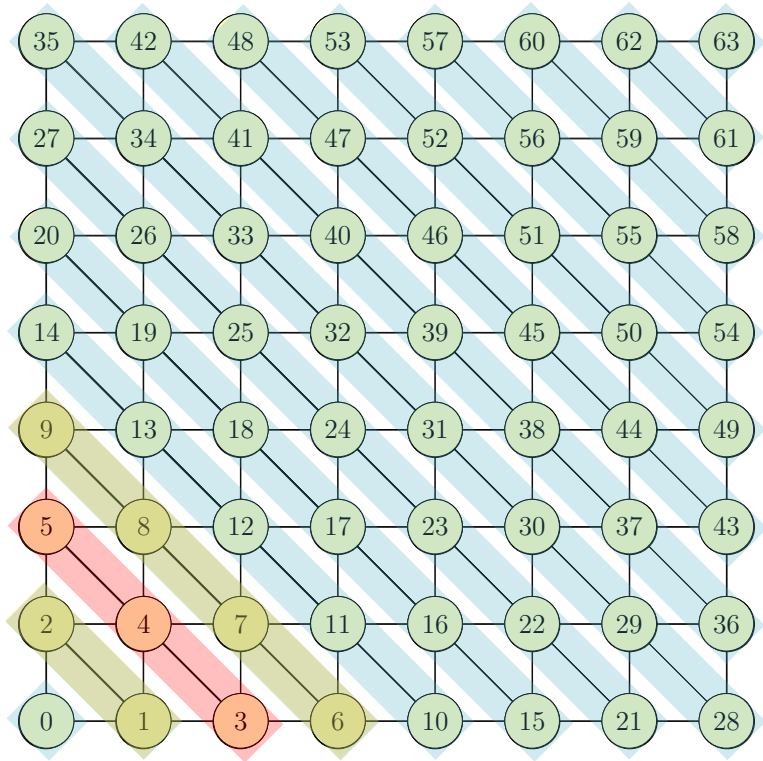
RACE approach



Matrix accessed 1 time from memory

How to do that in general for sparse matrices?

RACE



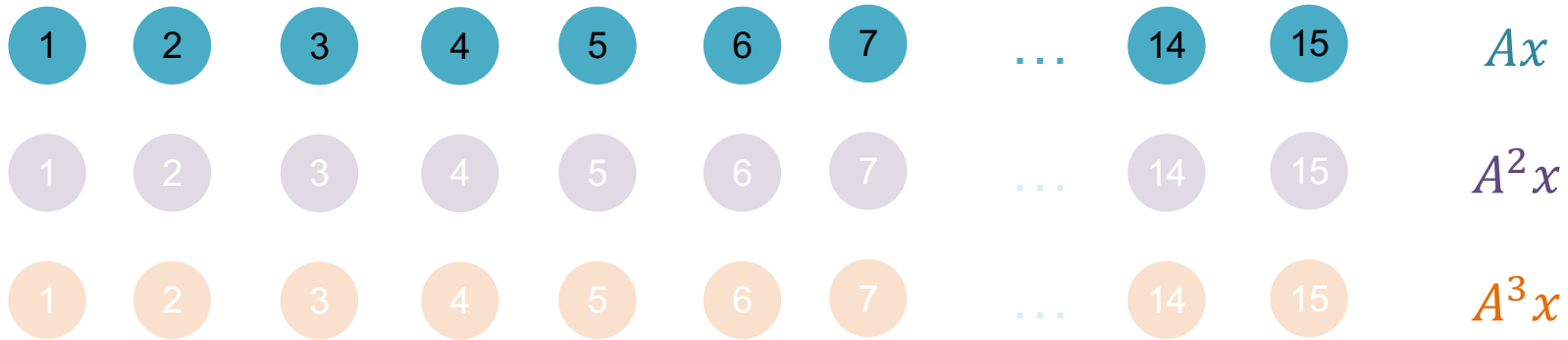
RACE – Level traversal and matrix powers

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

No cache blocking!

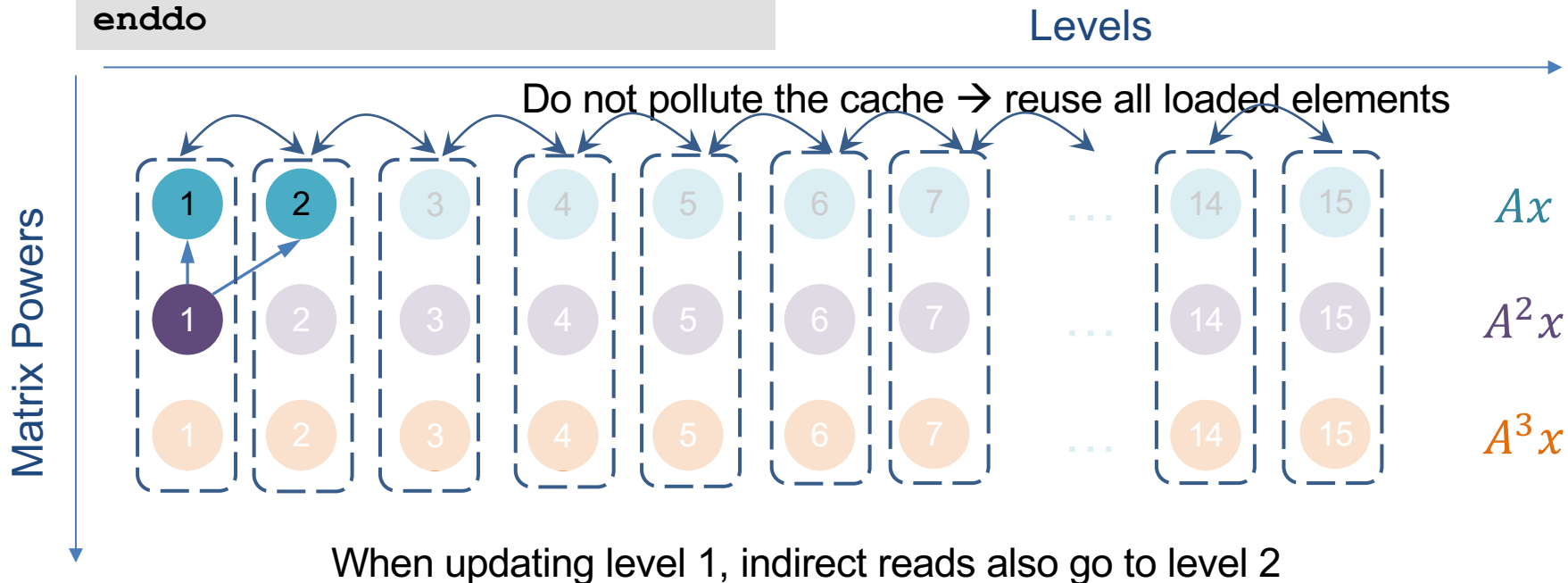
Levels

Matrix Powers



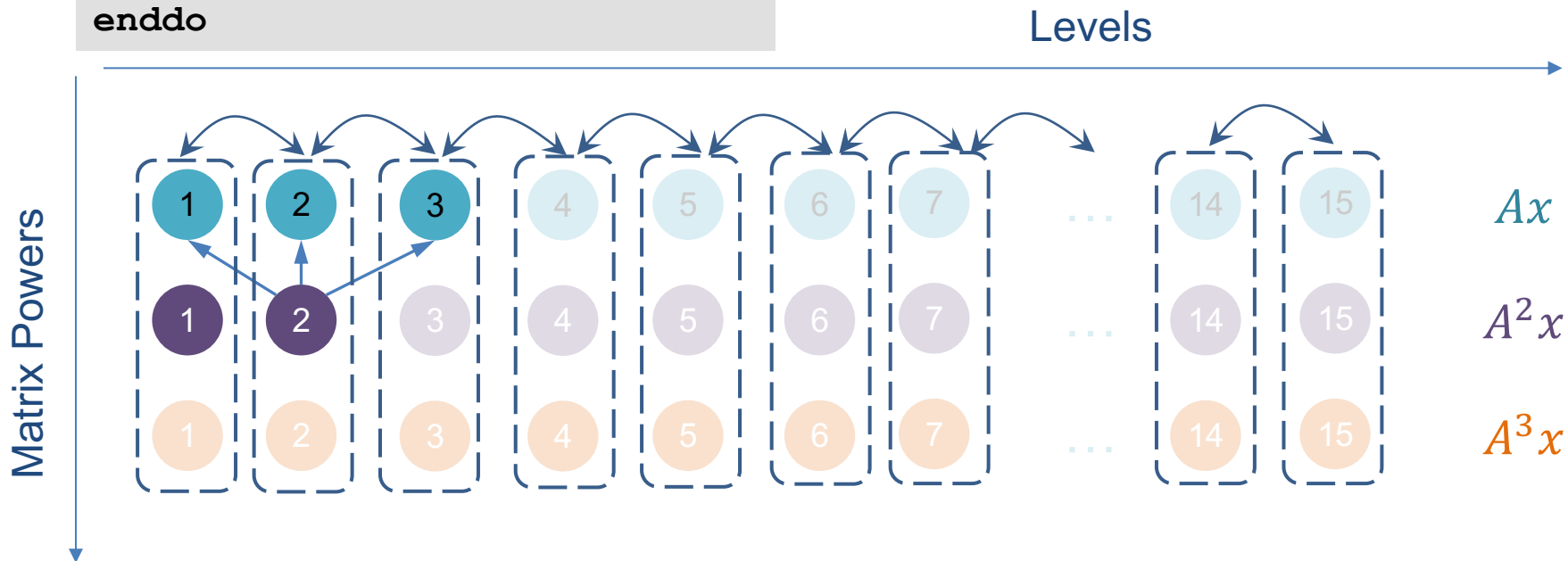
RACE – Level traversal and matrix powers

```
do k = 1, p
   $y(:, k) = \text{SpMV}(A, y(:, k-1))$ 
enddo
```



RACE – Level traversal and matrix powers

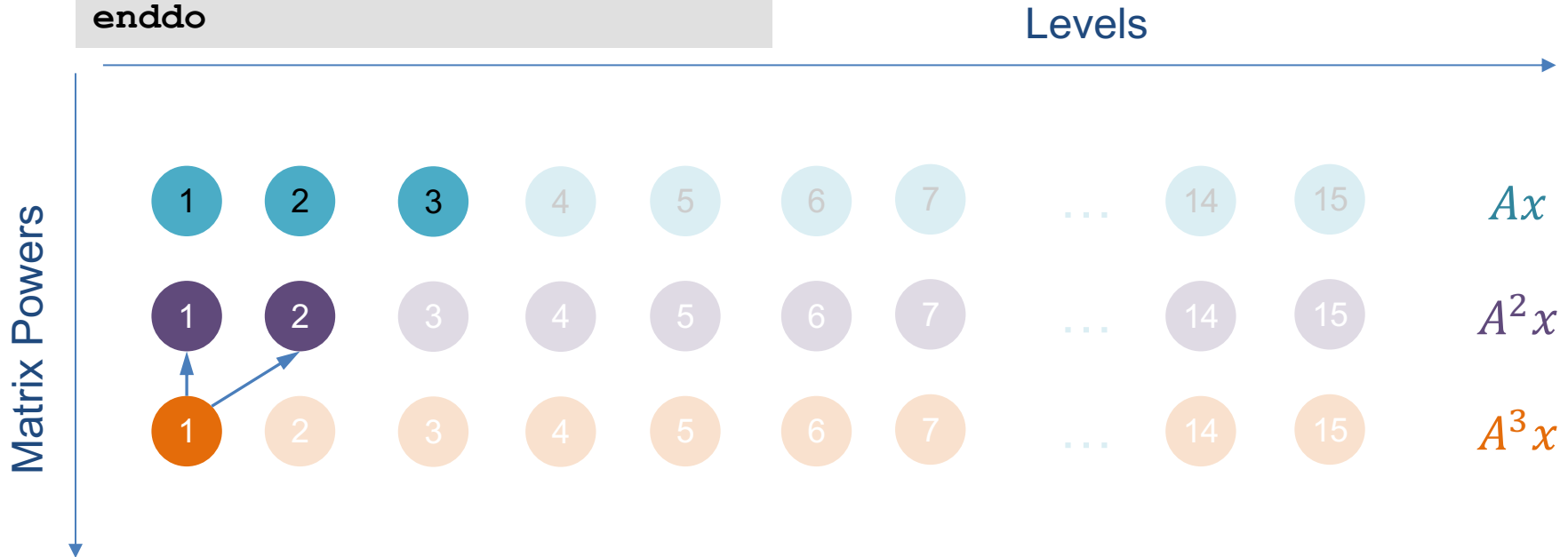
```
do k = 1, p  
   $\mathbf{y}(:, k) = \text{SpMV}(\mathbf{A}, \mathbf{y}(:, k-1))$   
enddo
```



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

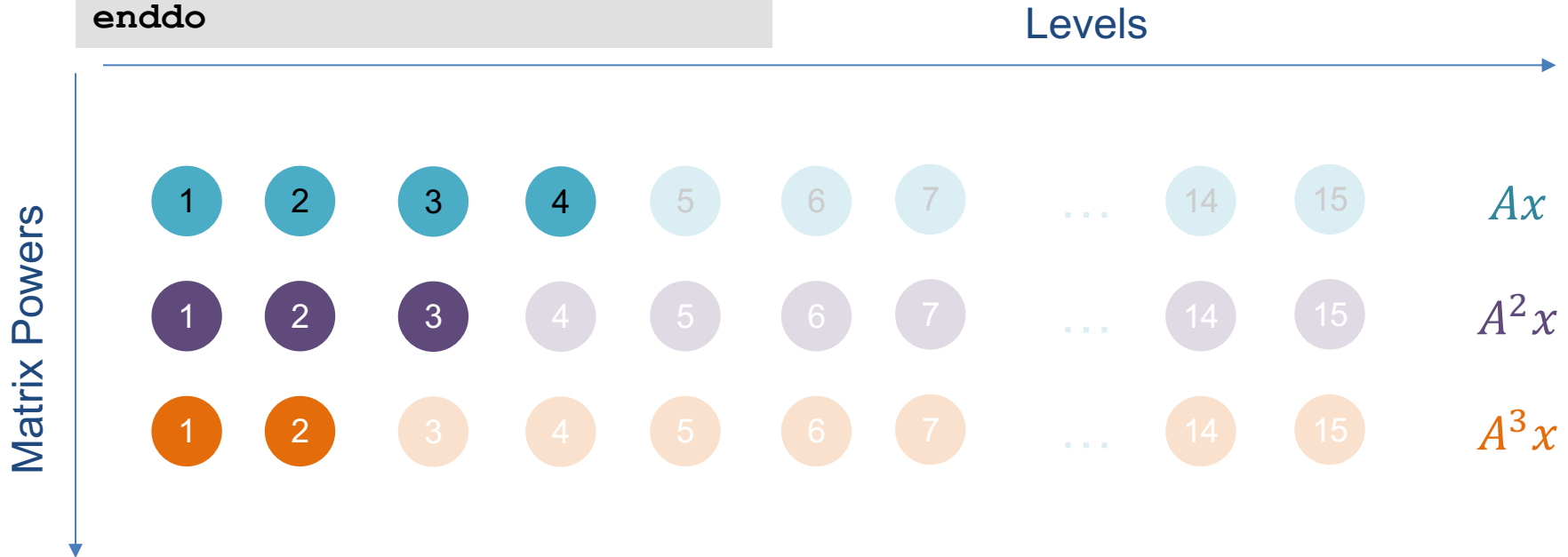
RACE – Level traversal and matrix powers

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



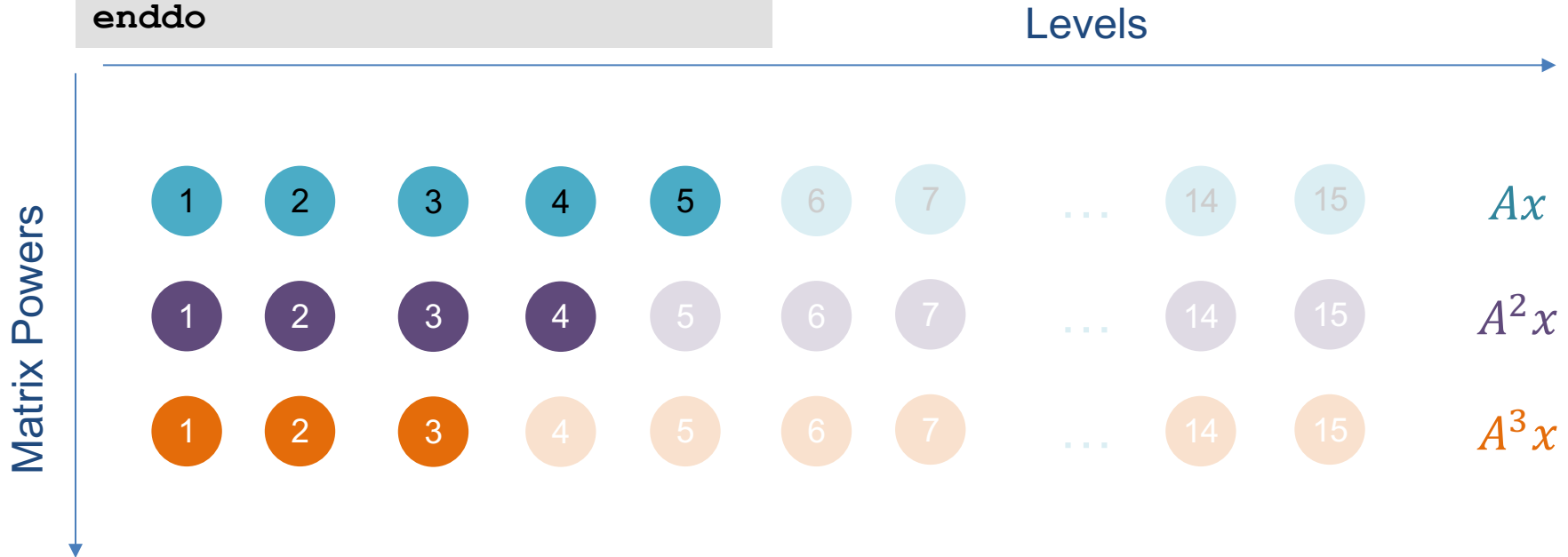
RACE – Level traversal and matrix powers

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



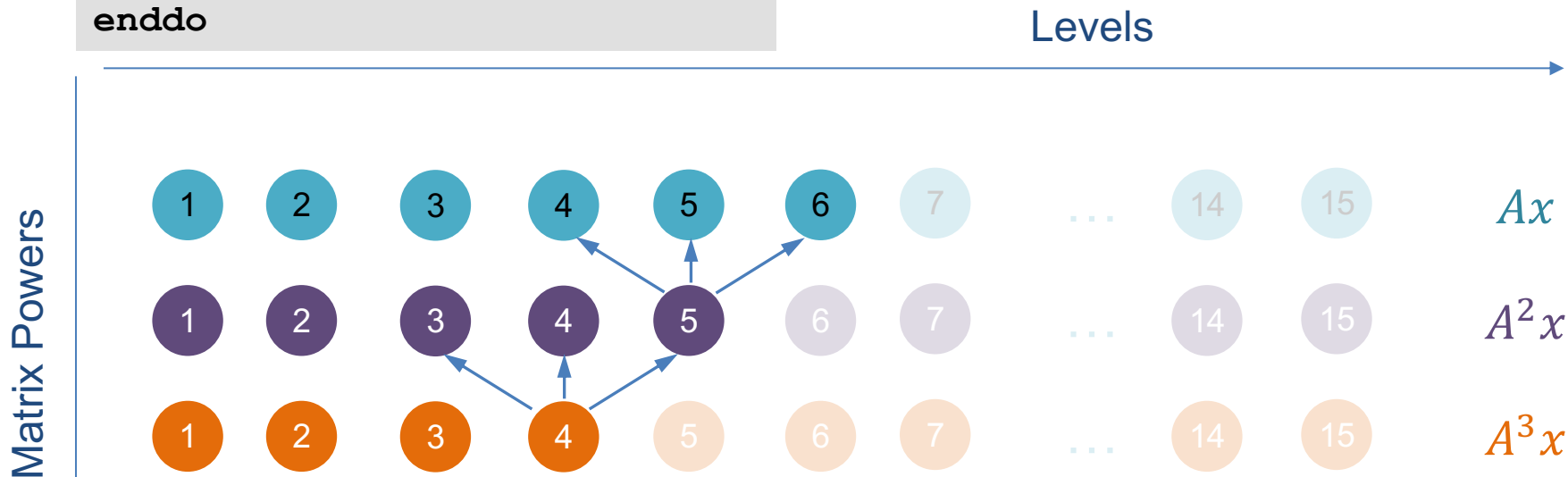
RACE – Level traversal and matrix powers

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



RACE – Level traversal and matrix powers

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



A is loaded only once if
 $(p + 1) \times N_{nz}(L) \times 12 \text{ bytes} < C$

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

RACE: MPK Pseudocode

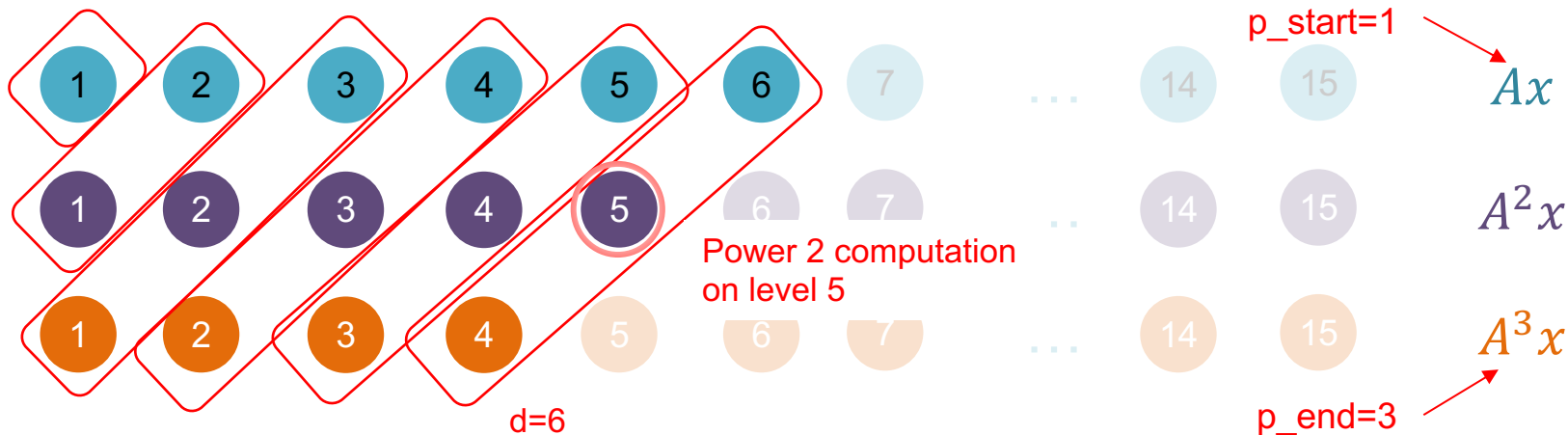
```
do d in 1:L+p-1
  p_start = max(1, d-(L-1))
  p_end = min(d, p)
  do k in p_start:p_end
    l=(d-k+1)
    y(:, k) = SpMV(A(j,:), y(:,k-1), level_ptr[l]:level_ptr[l+1])
  enddo
enddo
```

← Traverse along diagonal

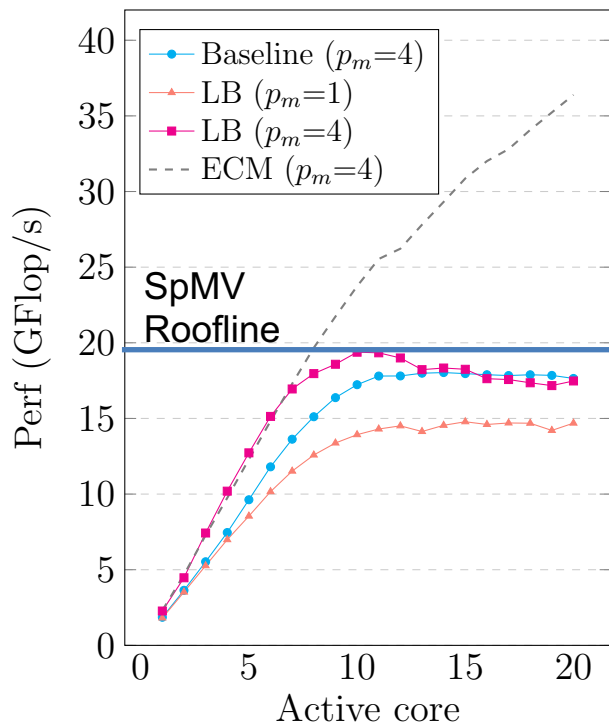
← All powers in diagonal

← Power k computation on level l

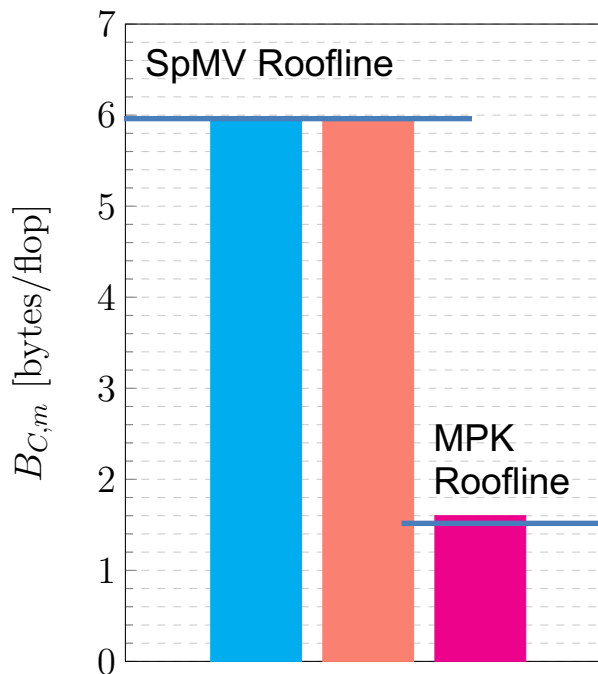
OpenMP parallel SpMV on all vertices in level l



RACE MPK – First Implementation



Performance



Memory traffic

Intel Xeon Gold 6248

- 1 Socket (20c)

pwtk matrix

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

RACE MPK – Performance Problem Identified

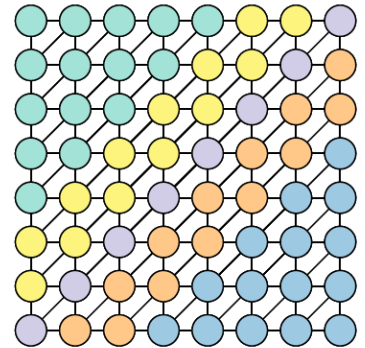
- Scheme seems to work (reduces data traffic) – at least for **pwtk**
- But: **Performance** ☹️ !!!!
- Analysis of hardware performance counters (LIKWID) for **pwtk** matrix:
`INSTR_RETIRED_ANY` up 2x for level based SpMV!

→ Frequent thread **synchronisations!**

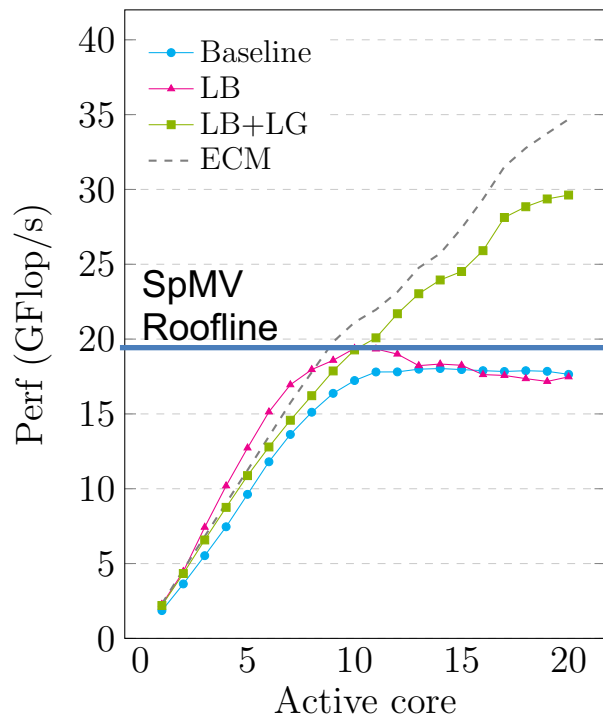
Reason: After each level threads sync!

Measures:

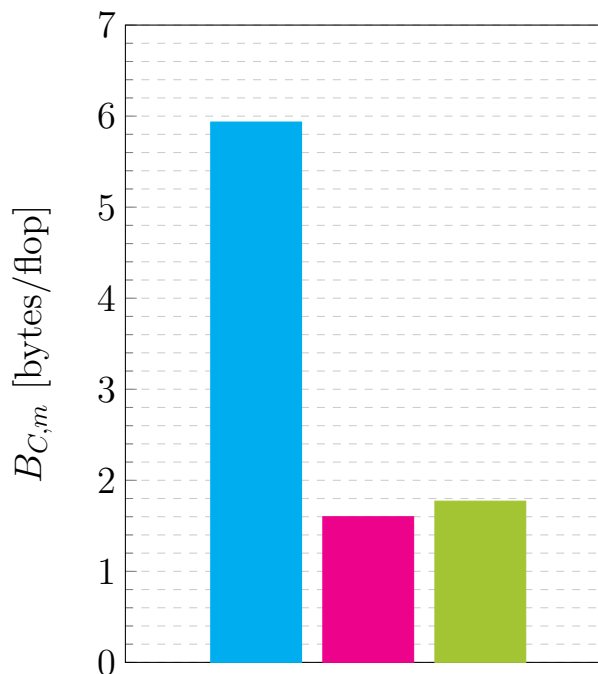
- Reduce #levels by level aggregation („**LG**“)
- Global sync. replaced by point-to-point sync. („**p2p**“)



RACE MPK – LG optimization



Performance



Memory traffic

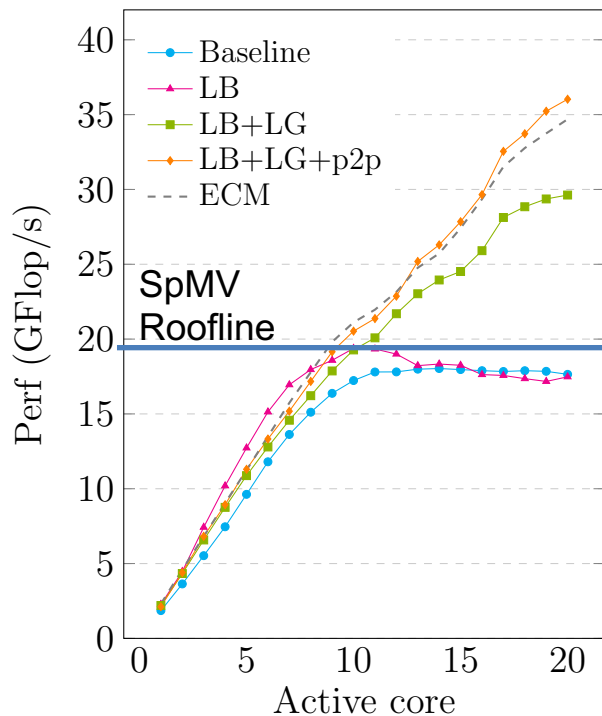
Intel Xeon Gold 6248

- 1 Socket (20c)

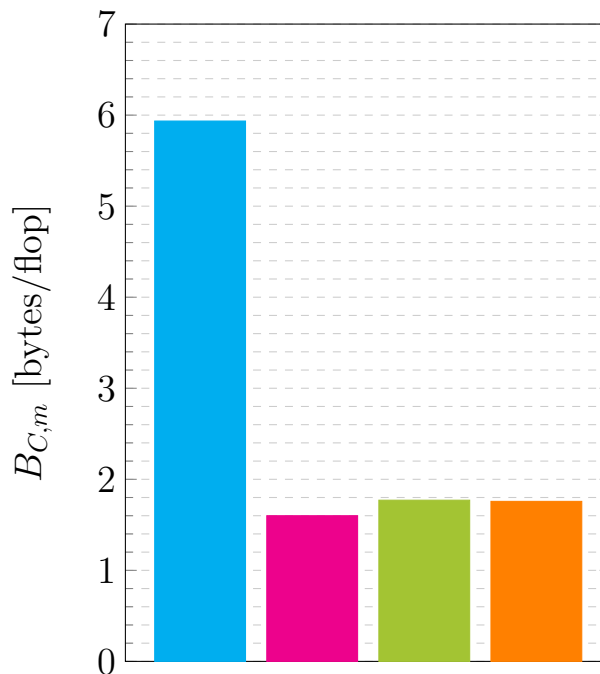
pwtk matrix

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

RACE MPK – LG+p2p optimization



Performance



Memory traffic

Intel Xeon Gold 6248

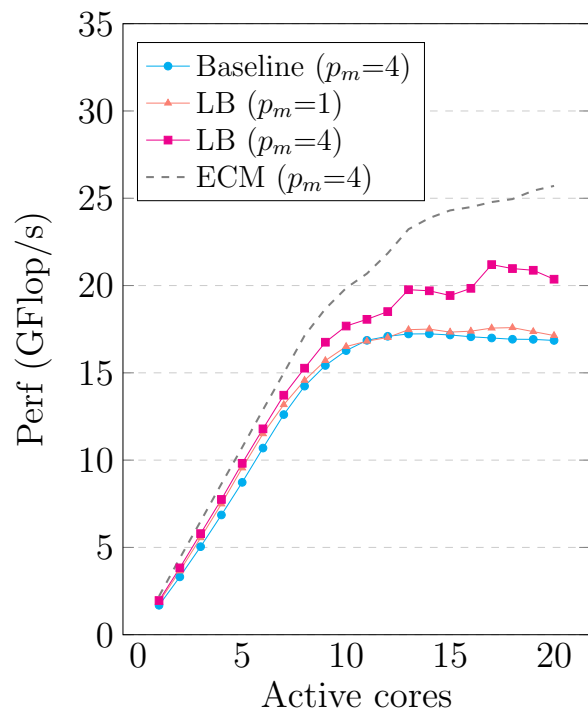
- 1 Socket (20c)

pwtk matrix

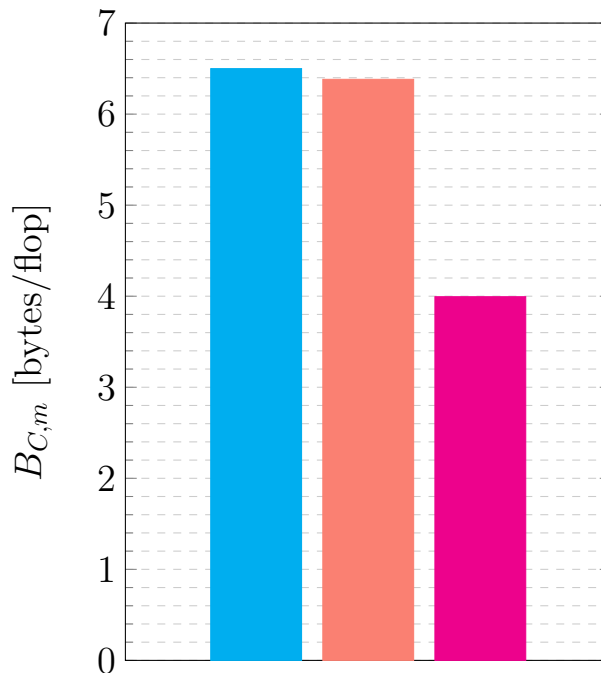
- $N_r = 217,918$
- $N_{nz} = 11,634,424$



RACE MPK – Yet another problem



Performance



Memory traffic

Intel Xeon Gold 6248

- 1 Socket (20c)

Flan_1565 matrix

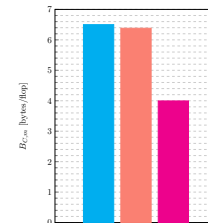
- $N_r = 1,564,794$
- $N_{nz} = 117,406,044$

Data traffic not reduced by factor of 4

Representative for large matrices!

RACE MPK – Performance Problem Identified (II)

- **Flan_1565** matrix – no significant adequate in data volume
- Analysis of level distribution for **Flan_1565** matrix:



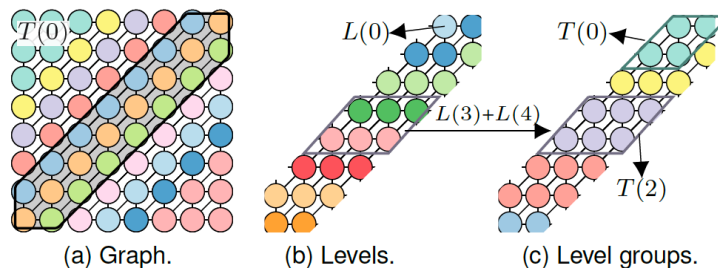
Few large levels → do not fit in cache!

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

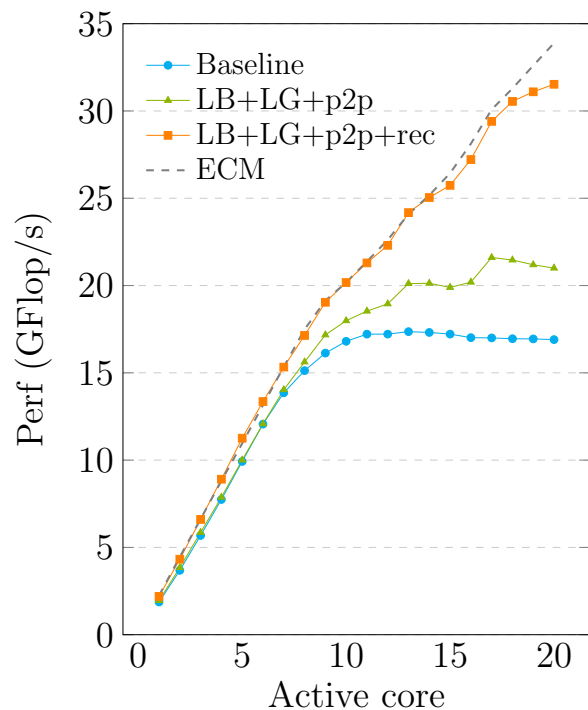
Counter - Measure:

Too big!!!

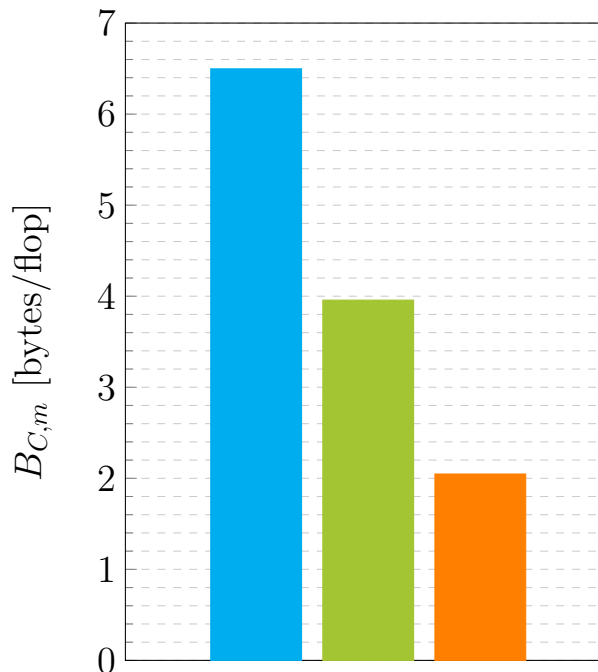
→ Apply RACE to a single / few levels recursively („rec“)!



RACE MPK – rec



Performance



Memory traffic

Intel Xeon Gold 6248

- 1 Socket (20c)

Flan_1565 matrix

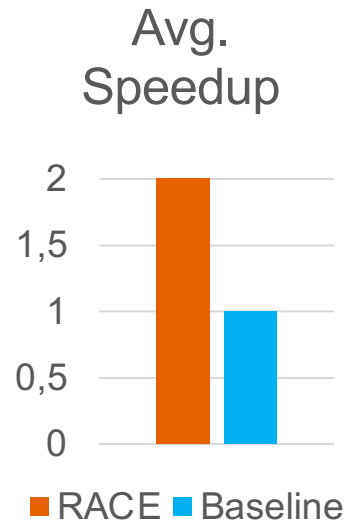
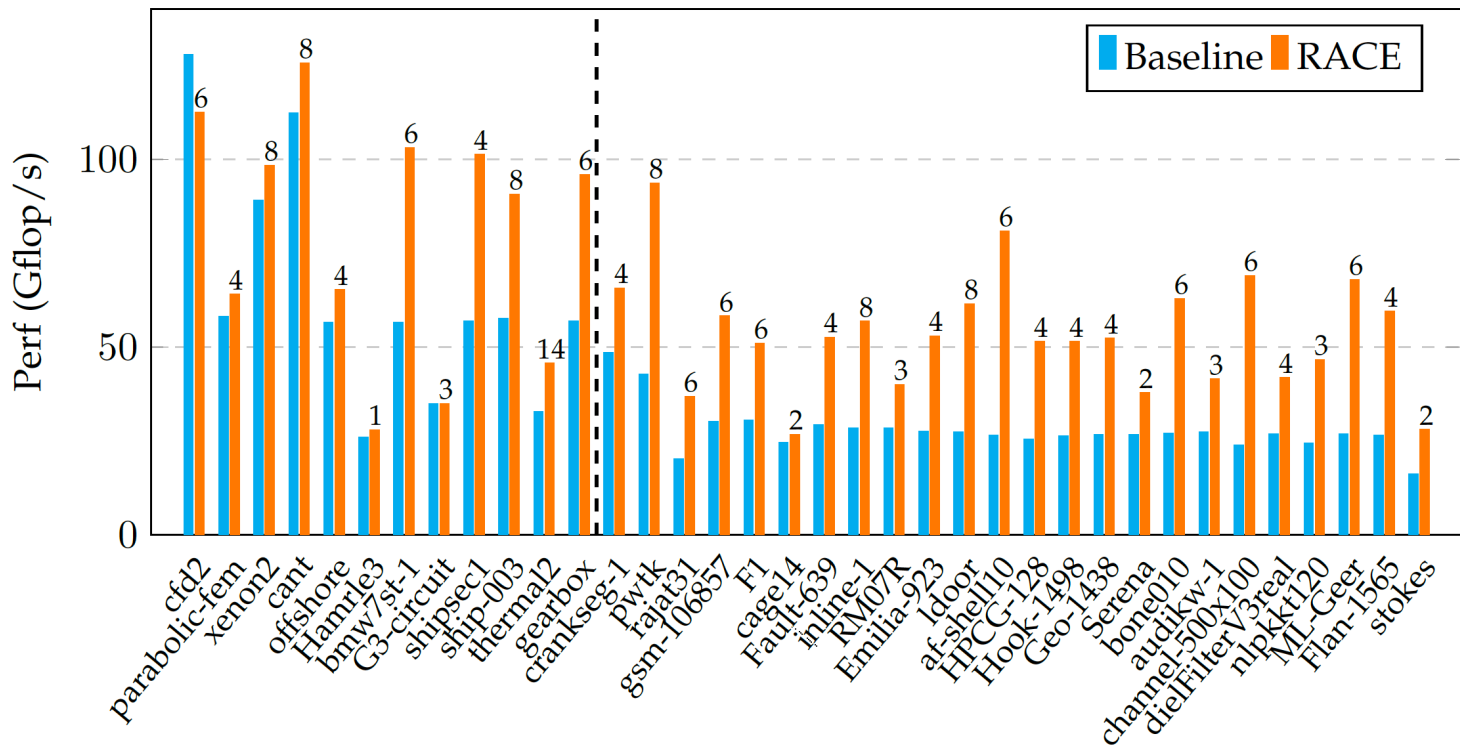
- $N_r = 1,564,794$
- $N_{nz} = 117,406,044$



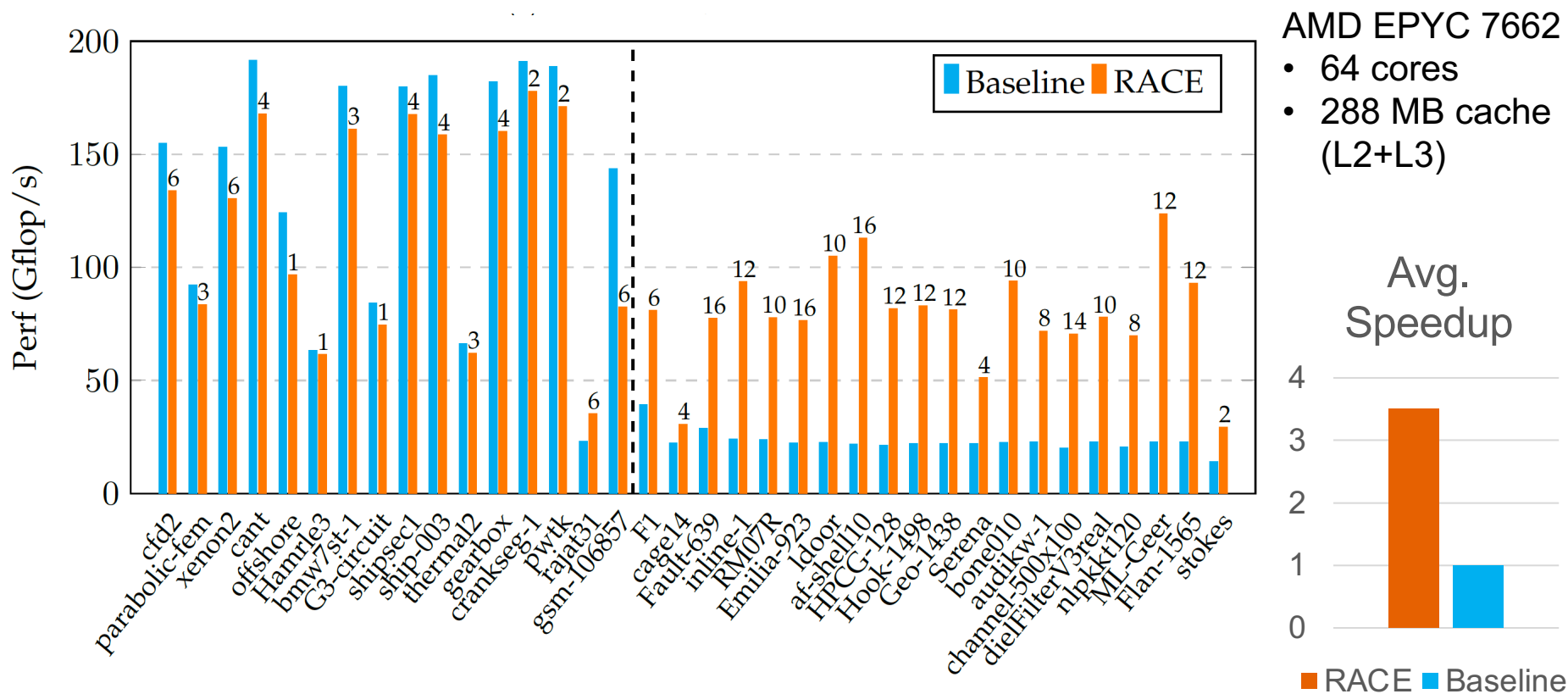
Matrix power kernel: Performance – Intel Ice Lake

Intel Xeon
Platinum 8368

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



Matrix power kernel: Performance – AMD Rome



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 \leftrightarrow 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
- Trilinos interface available

C. Alappat, J. Thies, G. Hager, H. Fehske, and G. Wellein: *Algebraic Temporal Blocking for Sparse Iterative Solvers on Multi-Core CPUs*. Submitted.
Preprint: [arXiv:2309.02228](https://arxiv.org/abs/2309.02228)