



LIBRSB: Multicore Sparse Matrix Performance across Languages and Architectures



Michele MARTONE

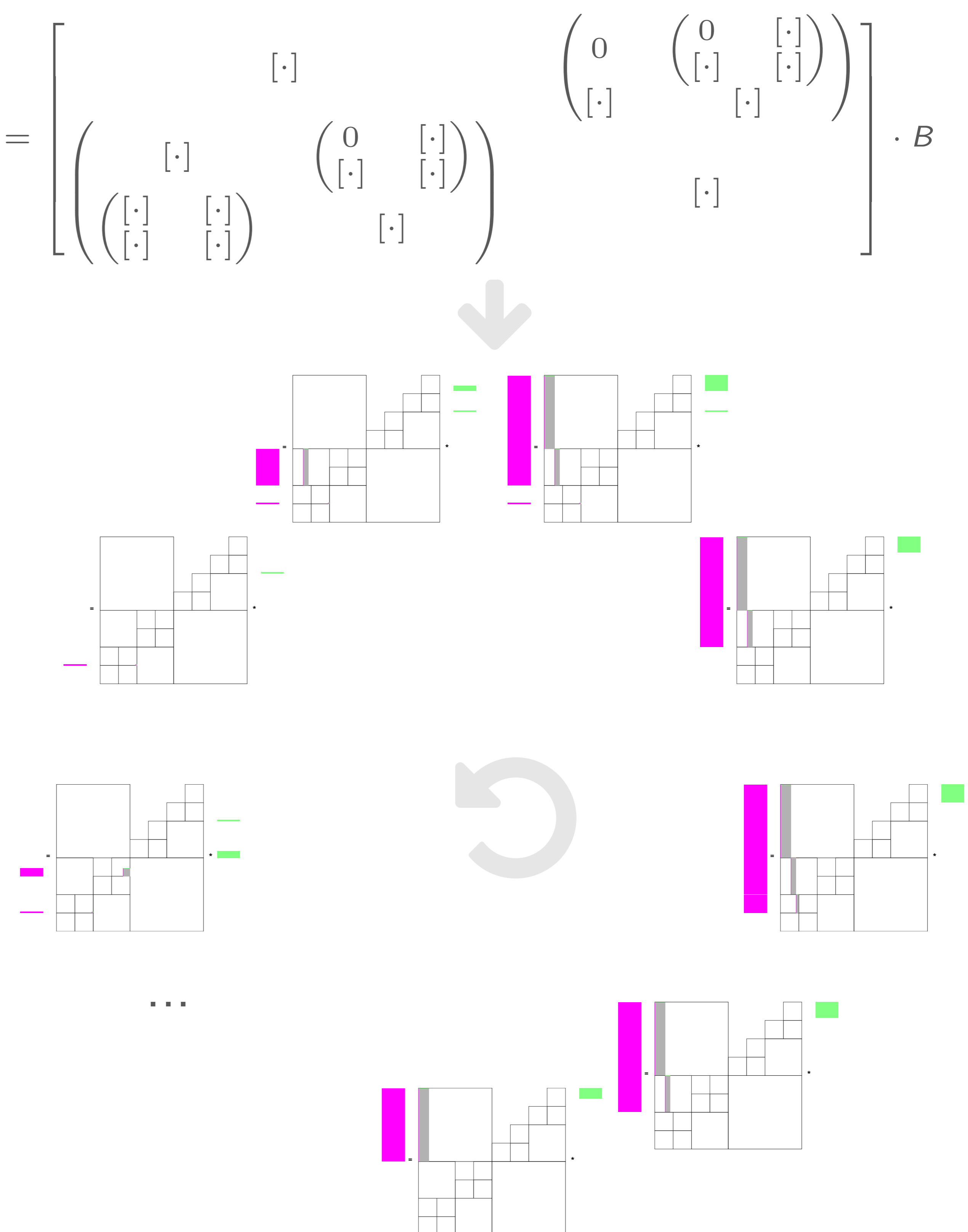
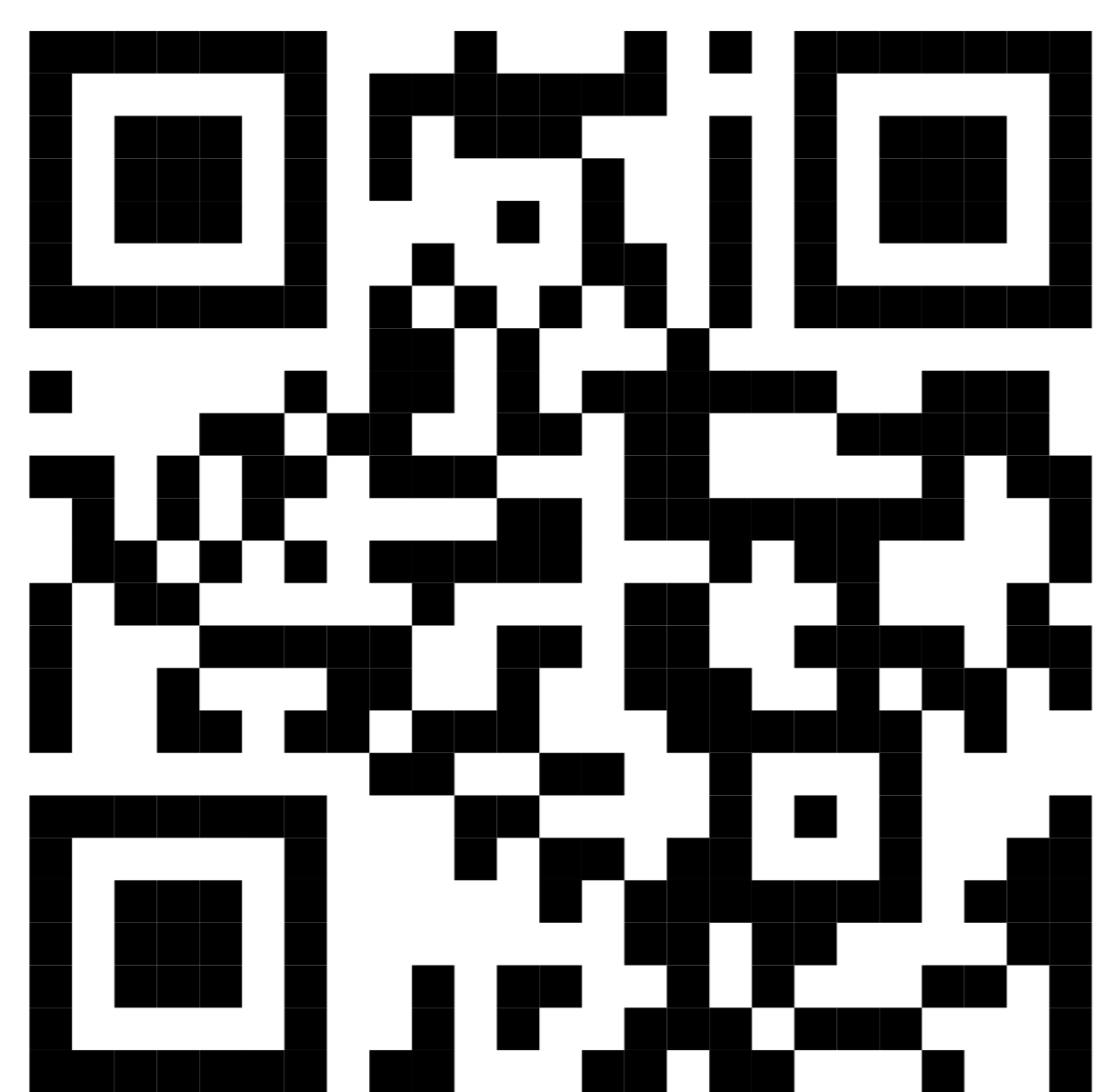
michele.martone@lrz.de

Leibniz Supercomputing Centre, Garching bei München, Germany

$$C = A \cdot B \rightarrow C = \begin{matrix} \text{sparse } A \\ \left[\begin{matrix} \cdot & \cdot \\ \cdot & \cdot \\ \cdot & \cdot \end{matrix} \right] \end{matrix} \cdot B \rightarrow C = \begin{bmatrix} \cdot & & \begin{pmatrix} 0 & \begin{bmatrix} \cdot \\ \cdot \end{bmatrix} \\ \cdot & \begin{bmatrix} \cdot \\ \cdot \end{bmatrix} \end{pmatrix} \\ \begin{pmatrix} \cdot & \cdot \\ \cdot & \cdot \end{pmatrix} & \begin{pmatrix} 0 & \cdot \\ \cdot & \cdot \end{pmatrix} & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix} \cdot B$$

LIBRSB: Universal Sparse BLAS Library

<https://librsb.sf.net>



apt install librsb-dev # spack install librsb # guix install librsb # eb search librsb #
Programming languages-related images with permission from www.octave.org, scipy.org, cplusplus.org, fortran-lang.org, wikipedia.org

A portable Sparse BLAS

“Sparse Basic Linear Algebra Subroutines”

for multicore CPUs (node-level)

one library, multiple APIs

this poster: LIBRSB overview and ongoing work

Numerical Techniques of Interest

iterative methods: *block Krylov*

require efficient Sparse Matrix-Matrix multiplication aka SpMM

SpMM in matrix form (*m* aka *NRHS* aka *number of right hand sides*):

$$\begin{bmatrix} c_{11} & \dots & c_{1m} \\ \vdots & \ddots & \vdots \\ c_{n1} & \dots & c_{nm} \end{bmatrix} \leftarrow \beta \begin{bmatrix} c_{11} & \dots & c_{1m} \\ \vdots & \ddots & \vdots \\ c_{n1} & \dots & c_{nm} \end{bmatrix} + \alpha \begin{bmatrix} a_{11} & \dots & a_{1m} \\ \vdots & \ddots & \vdots \\ a_{n1} & \dots & a_{nm} \end{bmatrix} \begin{bmatrix} b_{11} & \dots & b_{1m} \\ \vdots & \ddots & \vdots \\ b_{n1} & \dots & b_{nm} \end{bmatrix}$$

Modern C++ interface (rsb.hpp)

```
1 #include <rsb.hpp>
2 #include <vector>
3 #include <array>
4 using namespace ::rsb;
5
6 int main() {
7     RsbLib rsb;
8     const int nra { 7 }, nrhs { 2 };
9     const int nca { 6 }, ncb { 6 };
10    const std::vector<int> IA {0,1,2,3,4,5,1};
11    const int JA [ ] = {0,1,2,3,4,5,0};
12    const std::vector<double> VA {1,1,1,1,1,1,2}, B(nrhs * nca,1);
13    std::array<double,nrhs * nra> C;
14    const double alpha {2}, beta {1};
15
16    RsbMatrix<double> mtx(IA,JA,VA,nra); // Notice above declarations of IA,JA,VA
17
18    mtx.spmm(RSB_TRANSPOSITION_N, alpha, nrhs, RSB_FLAG_WANT_ROW_MAJOR_ORDER, B,
19            beta, C); // ditto for B and C
20 }
```

Figure 2: RsbMatrix's spmm() and others use C++20's std::span. Gives freedom in vectors' type choice. Wraps rsb.h.

New: ongoing work

Improve performance of SpMM CSR kernels:

AVX512 intrinsics (*by-rows*, unsymmetric, untransposed)

more specializations

autotuning challenges: e.g. use CSR if more performant than RSB

LIBRSB

project page: <http://librsb.sf.net>

>100KLOC of C99, OPENMP, and modern templated C++

node-level **shared-memory**-parallel operations, for:

Sparse BLAS: matrix assembly/destroy, SpMM, triangular solve

interactive applications (update, sparse-sparse ops, conversions)

distributed-memory applications (block extract, update, etc.)

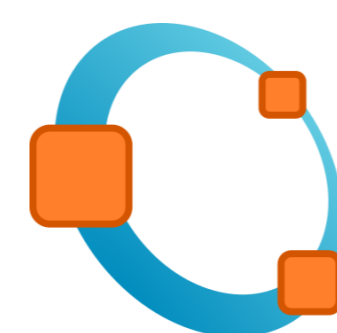
with own interface in C/C++ and FORTRAN

with Sparse BLAS interface (BLAS Technical Forum Standard)

with interfaces for GNU OCTAVE and PYTHON interpreters

LGPLv3-licensed free software, available via SPACK, GUIX-HPC, EasyBuild, and on Debian, Ubuntu, OpenSUSE, Windows,...

GNU Octave + liboctave + LIBRSB = SparseRSB



GNU OCTAVE: a MATLAB-like interactive numerical language

liboctave: access OCTAVE via C++

```
1 octave:1> A = spsersb ( sparse ( rand (3) > .6) )
2 A =
3
4 Recursive Sparse Blocks (rows = 3, cols = 3, nnz = 2 [22%])
5
6 (3, 1) -> 1
7 (3, 3) -> 1
```

Figure 3: The spsersb usage is styled after the sparse built-in. So most of operators (*,*,=(,(:),...) work the same way.

Recursive Sparse Blocks (RSB) Layout

for large matrices (uses cache locality, coarse thread parallelism)

supports autotuning (layout adjusted to maximize performance)

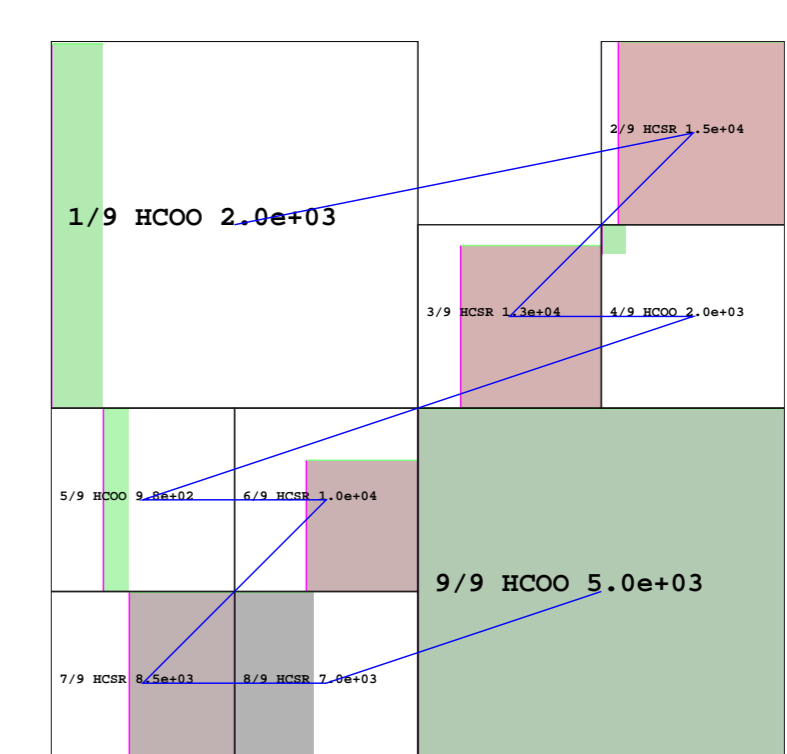


Figure 1: Instance of classical test matrix *bayer02* ($14k \times 14k$, $64k$ nonzeros). Black-bordered boxes are *sparse blocks*, and are *Z-ordered*. Greener have fewer nnz than average, redder have more. Either in “Coordinate” format (COO) or “Compressed Sparse Rows” (CSR). Blocks rows (LHS) and columns (RHS) ranges evidenced (left and top side).

Python + Cython + LIBRSB = PyRSB



SciPy: popular PYTHON scientific computing API

CYTHON: optimising static compiler for C extensions to PYTHON

```
1 import numpy
2 import scipy
3 from rsb import rsb_matrix
4 V=[ 11.,12.,22.]
5 I=[ 0, 0, 1]
6 J=[ 0, 1, 1]
7 a = rsb_matrix((V, (I, J)), [3,3])
8 ...
9 y = y + a * x;
```

Figure 4: API styled after SciPy's widely known csr_matrix. Additionally, it offers RSB-specific functionality, e.g. blocks autotuning.

Acknowledgements

This work has received funding by PRACE IP6/WP8 and PRACE High Level Support Teams.

Development and test and of LIBRSB utilizes systems of the BEAST (Bavarian Energy Architecture & Software Testbed) at LRZ.

