

# A Unified Sparse Matrix Storage Format for Heterogeneous Systems

Moritz Kreutzer, Georg Hager, and Gerhard Wellein

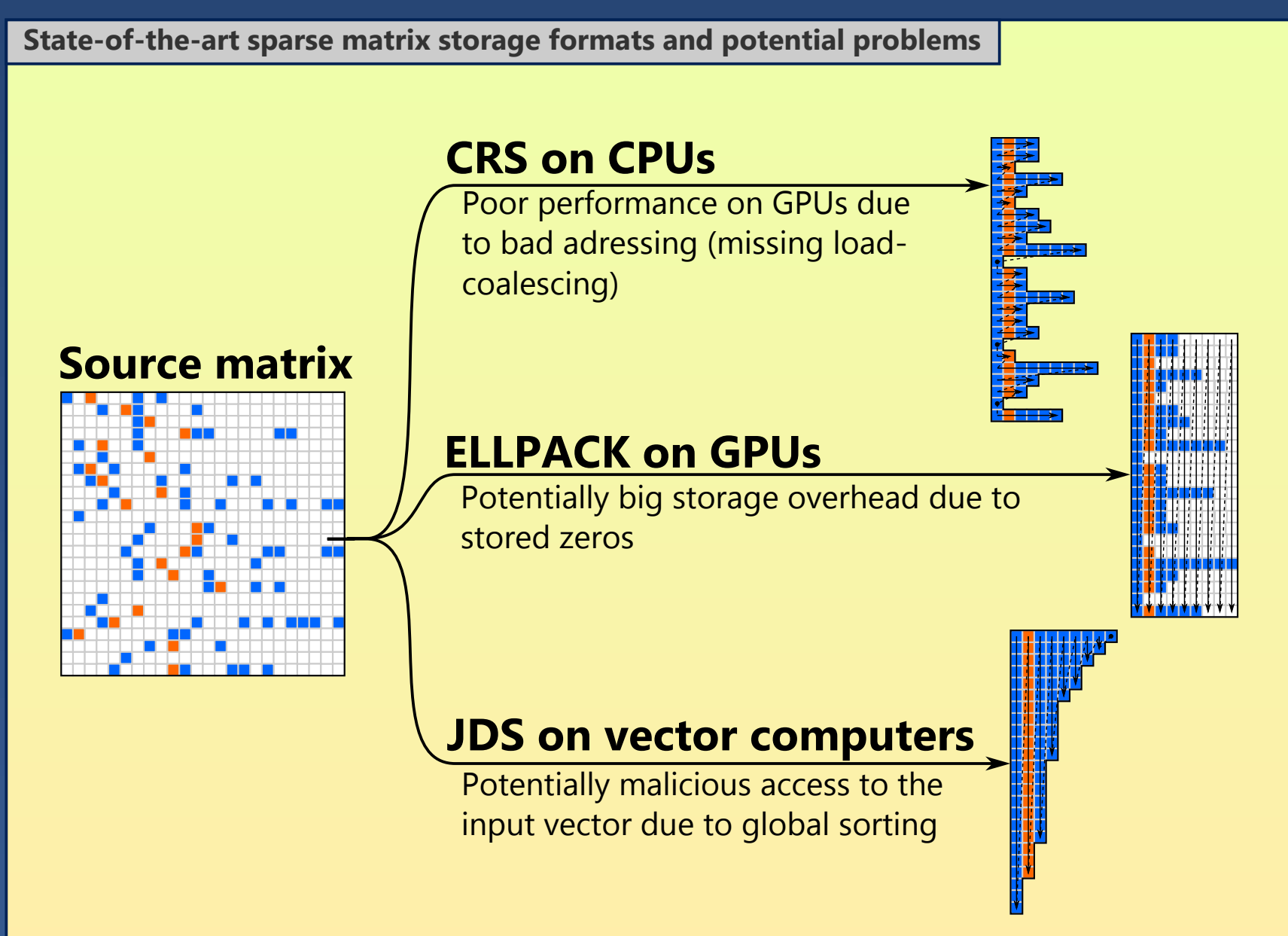
{moritz.kreutzer,georg.hager,gerhard.wellein}@fau.de

Erlangen Regional Computing Center, Friedrich-Alexander University of Erlangen-Nuremberg, Germany

## Introduction

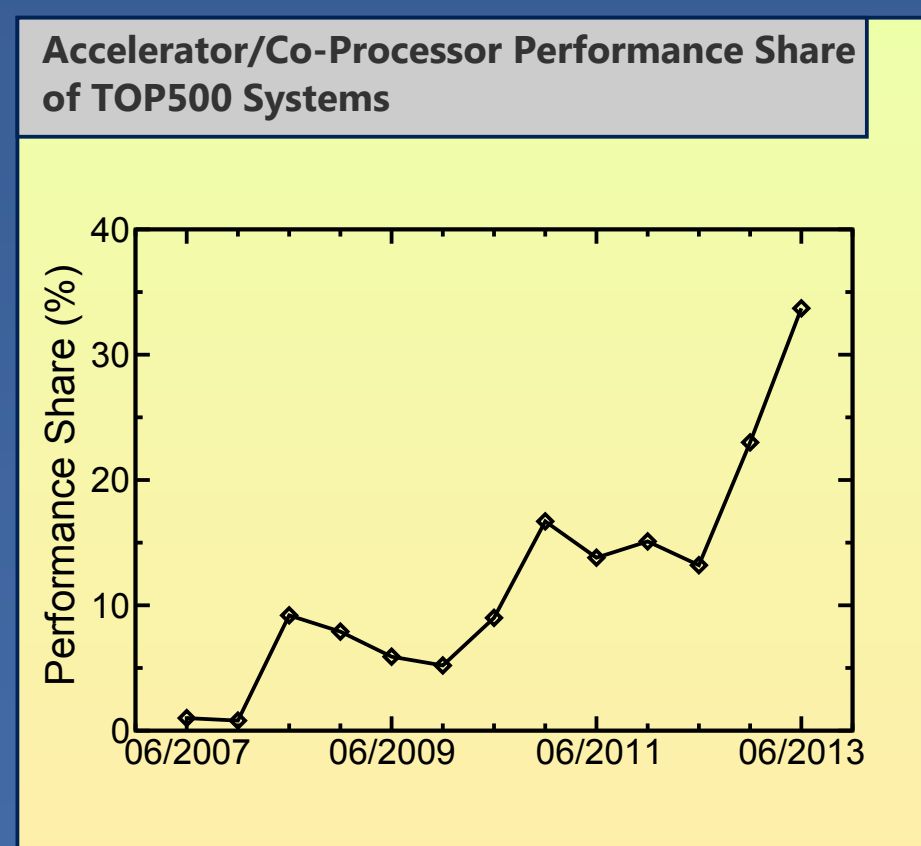
Sparse Matrix Vector Multiplication (SpMVM,  $y=Ax+y$ ) is a very bandwidth-demanding operation and one of the most time-consuming routines in many scientific applications. Sparse matrices (number of non-zeros  $N_{nz}=O(N)$ ) are usually stored in a way which omits (most of) the zero entries in order to minimize the storage overhead.

The SpMVM performance relies on the optimal storage format which is closely related to the compute architecture in use. The following image illustrates the state-of-the-art storage formats for three different architectures and potential problems that may arise for them.



Emerging co-processors/accelerators like the Intel Xeon Phi or Nvidia Tesla GPUs are of special interest for the spMVM because of their large memory bandwidth together with a high level of on-chip parallelism.

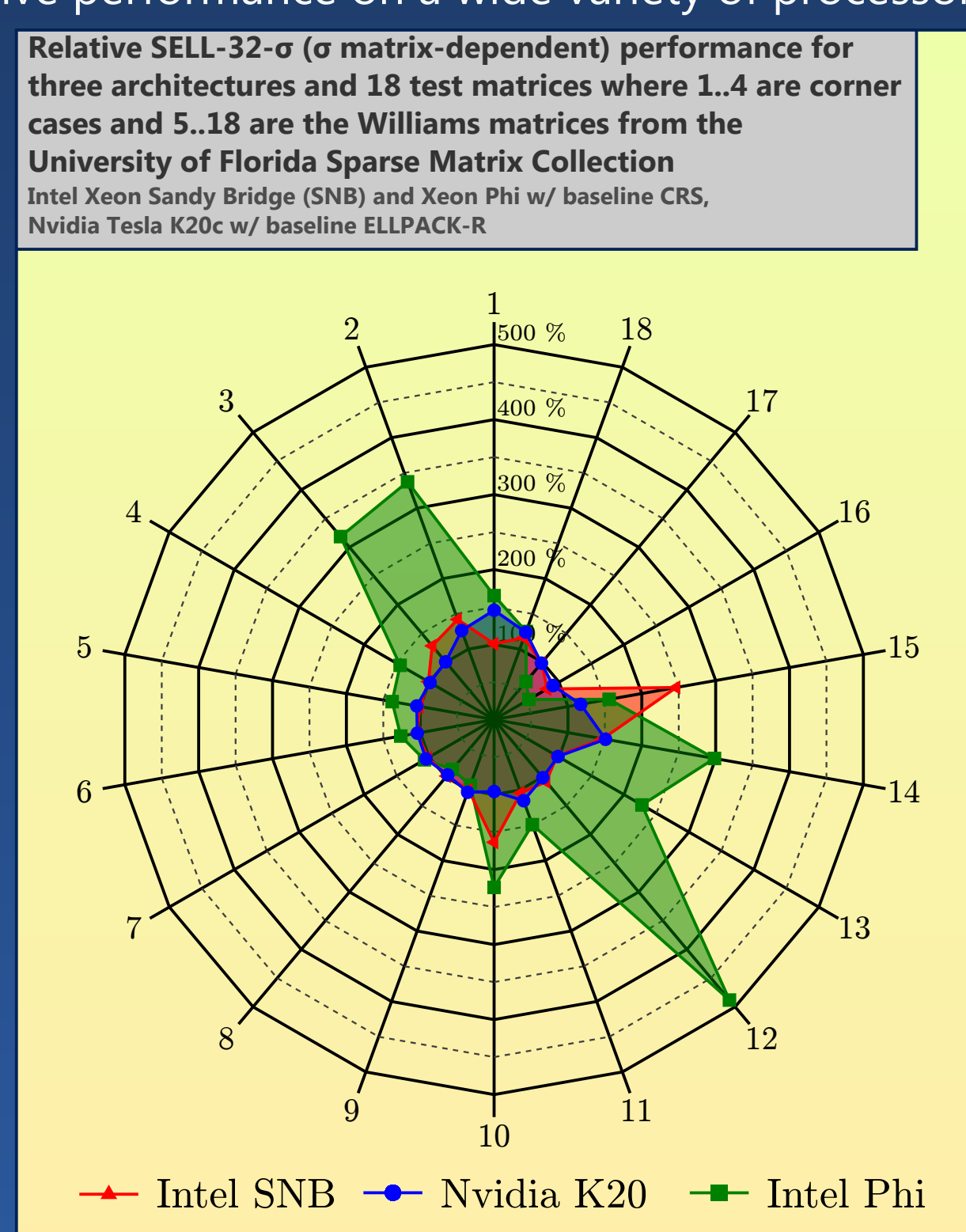
Upcoming cluster systems will probably be more and more of heterogeneous nature, i.e., consisting of standard CPUs together with accelerators. Sustainable and modern high-performance software should be able to utilize all parts of a heterogeneous cluster. Hence, establishing a single storage format for sparse matrices that yields good performance on all architectures is of broad interest.



## Contribution

This work demonstrates the feasibility of a unified storage format for sparse matrices, namely **SELL-C- $\sigma$** . It builds on Sliced ELLPACK, which has only been used on GPUs up to now, and provides competitive performance on a wide variety of processor designs.

It has two parameters,  $C$  and  $\sigma$ , for which good values can be predicted in advance. A performance model allows us to study the interplay of  $C$ ,  $\sigma$  and further matrix properties such as the "chunk occupancy" and the number of non-zeros per row.



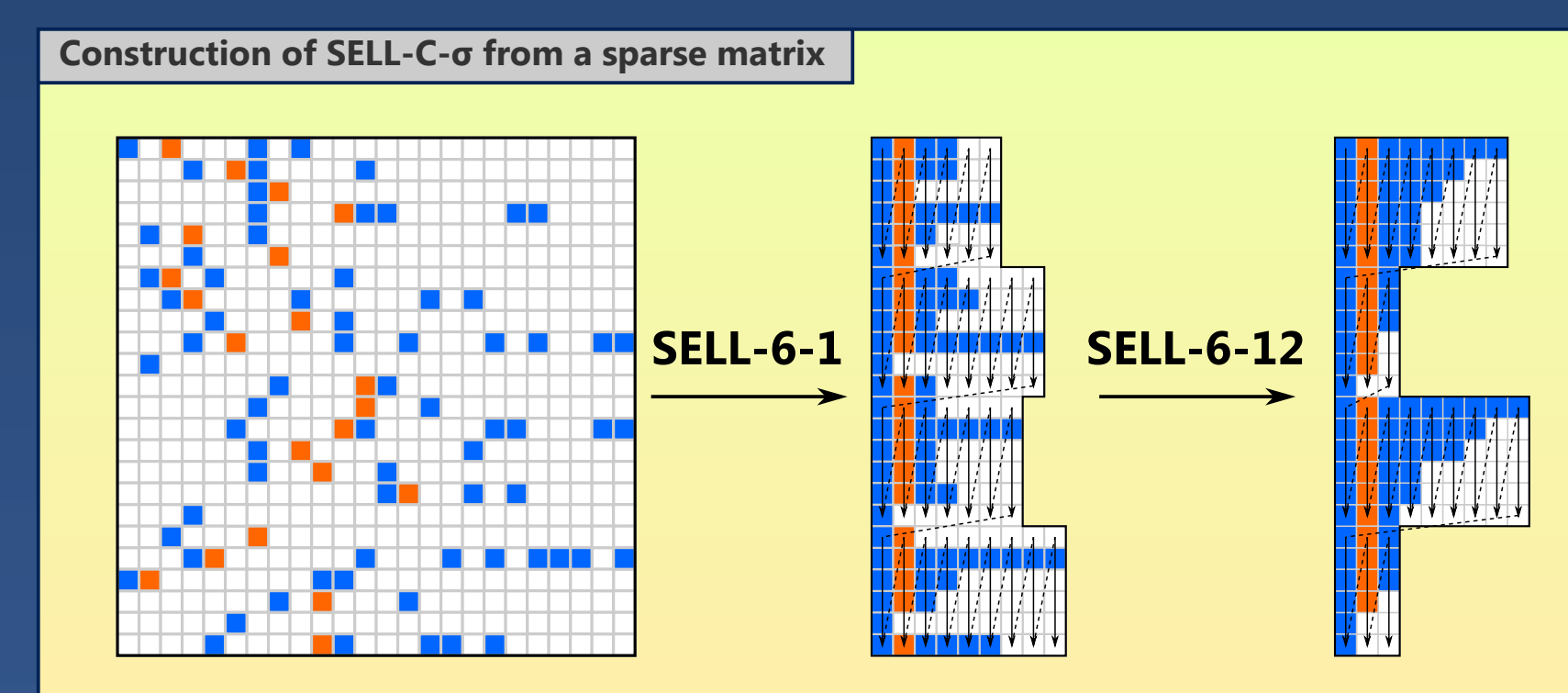
## SELL-C- $\sigma$ construction

Non-zero matrix entries and column indices are stored column-wise in chunks of height  $C$  ( $C=1$ : CRS). All rows of a chunk are padded with zeros to have the same length. Storage overhead is minimal and performance is usually best if the chunk structure is chosen in accordance with the width of a SIMD (Single Instruction Multiple Data) register on x86 and Intel MIC.

On GPUs, SIMD $\Delta$ SIMT (Single Instruction Multiple Threads) with width = warp size (32 threads on Nvidia Kepler). In a heterogeneous setting: Choose  $C$  in accordance with the largest available SIMD width.

Architecture	SIMD width
Intel Sandy Bridge	256 bits
Intel Xeon Phi	512 bits
Nvidia Kepler	2048 bits

Sorting rows by the number of non-zero entries within a limited scope  $\sigma$  of rows reduces the overhead as equally large rows get close to each other and can increase performance. However, if  $\sigma$  is too large (extremal case: global sorting) performance may drop (potentially malicious access pattern to  $x$ ).



SELL-C- $\sigma$  data (besides  $C$  and  $\sigma$ ):

- `double val[]` matrix values
- `double col[]` column indices
- `int chunkWidth[]` length of the chunk's longest row
- `int chunkStart[]` each chunk's starting offset in val/col
- `int nRows` number of rows (multiple of  $C$ )

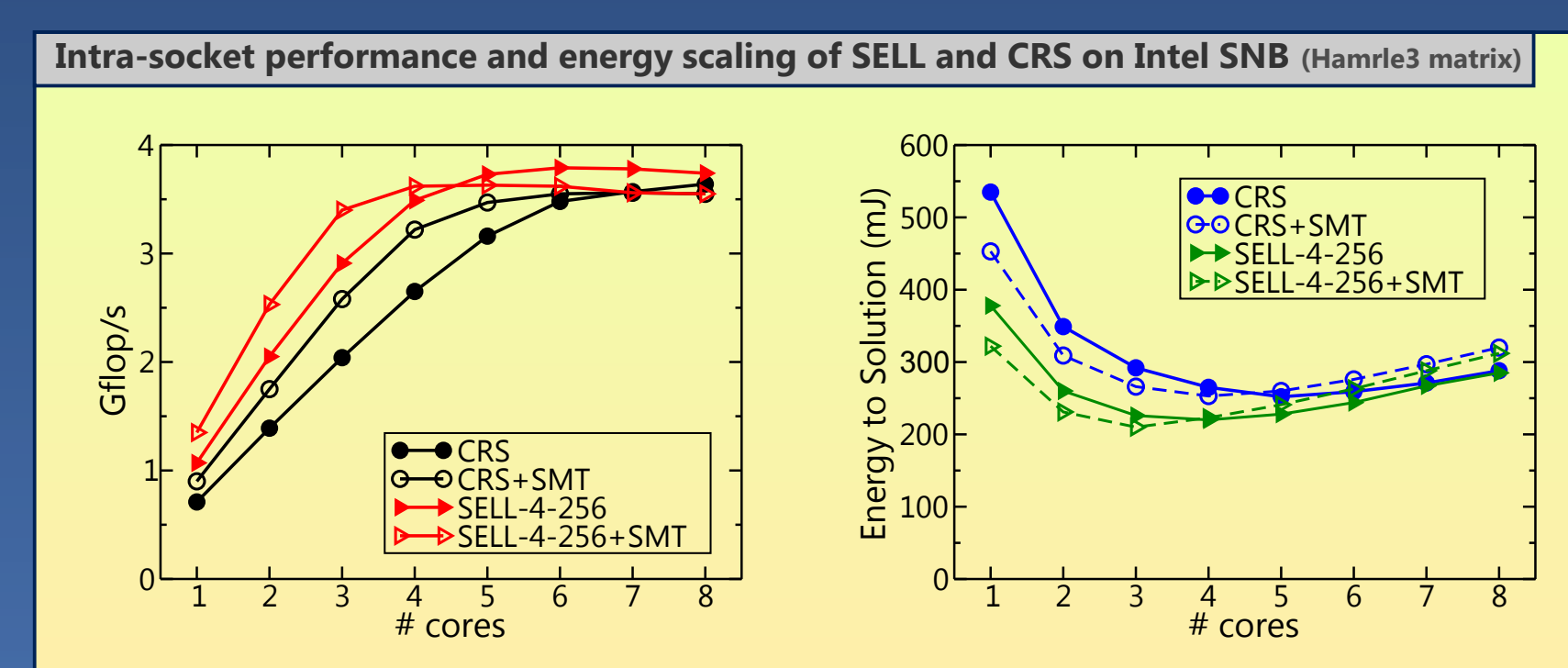
The SELL-C- $\sigma$  spMVM kernel can be easily vectorized by the compiler and requires neither a remainder loop nor a vector reduction, in contrast to the vectorized CRS kernel.

```

SELL-4- $\sigma$  kernel in C
for (i = 0; i < nRows/4; ++i)
{ // loop over chunks
  for (j = 0; j < chunkWidth[i]; ++j)
  { // loop inside chunk
    y[i*4+0] += val[chunkStart[i]+j*4+0] *
                x[col[chunkStart[i]+j*4+0]];
    y[i*4+1] += val[chunkStart[i]+j*4+1] *
                x[col[chunkStart[i]+j*4+1]];
    y[i*4+2] += val[chunkStart[i]+j*4+2] *
                x[col[chunkStart[i]+j*4+2]];
    y[i*4+3] += val[chunkStart[i]+j*4+3] *
                x[col[chunkStart[i]+j*4+3]];
  }
}
    
```

## Scaling and Energy

SpMVM is a bandwidth-bound kernel for large data sets and obviously, SELL-C- $\sigma$ 's SIMD capabilities cannot increase the bandwidth and saturated spMVM performance. However, due to faster in-socket scaling a considerable amount of energy can be saved in SELL-C- $\sigma$  compared to CRS. Using Simultaneous Multi-Threading (SMT) can further improve the energy efficiency.



## Performance Model

The "chunk occupancy"  $\beta$  quantifies the overhead by the zero padding. It is the ratio between  $N_{nz}$  and the elements stored in the SELL-C- $\sigma$  format ( $N_C$  denotes the number of chunks):

$$\beta = \frac{N_{nz}}{\sum_{i=0}^{N_C} C \cdot \text{chunkWidth}[i]}$$

$$\beta_{\text{worst}} = \frac{\sum_{i=0}^{N_C} (N+C-1)}{\sum_{i=0}^{N_C} CN} = \frac{N+C-1}{CN} \gg \frac{1}{C} \quad \text{matrix with one line of length } N \text{ and } C-1 \text{ lines of length 1 per chunk}$$

$$\beta_{\text{best}} = \frac{\sum_{i=0}^{N_C} C \cdot \text{chunkWidth}[i]}{\sum_{i=0}^{N_C} C \cdot \text{chunkWidth}[i]} = 1 \quad \text{matrix with all lines of equal length in each chunk}$$

The code balance of the SELL spMVM kernel for large data sets can be given as follows (8-byte values, 4-byte indices,  $N_{nz}$  denoting the number of non-zeros per row):

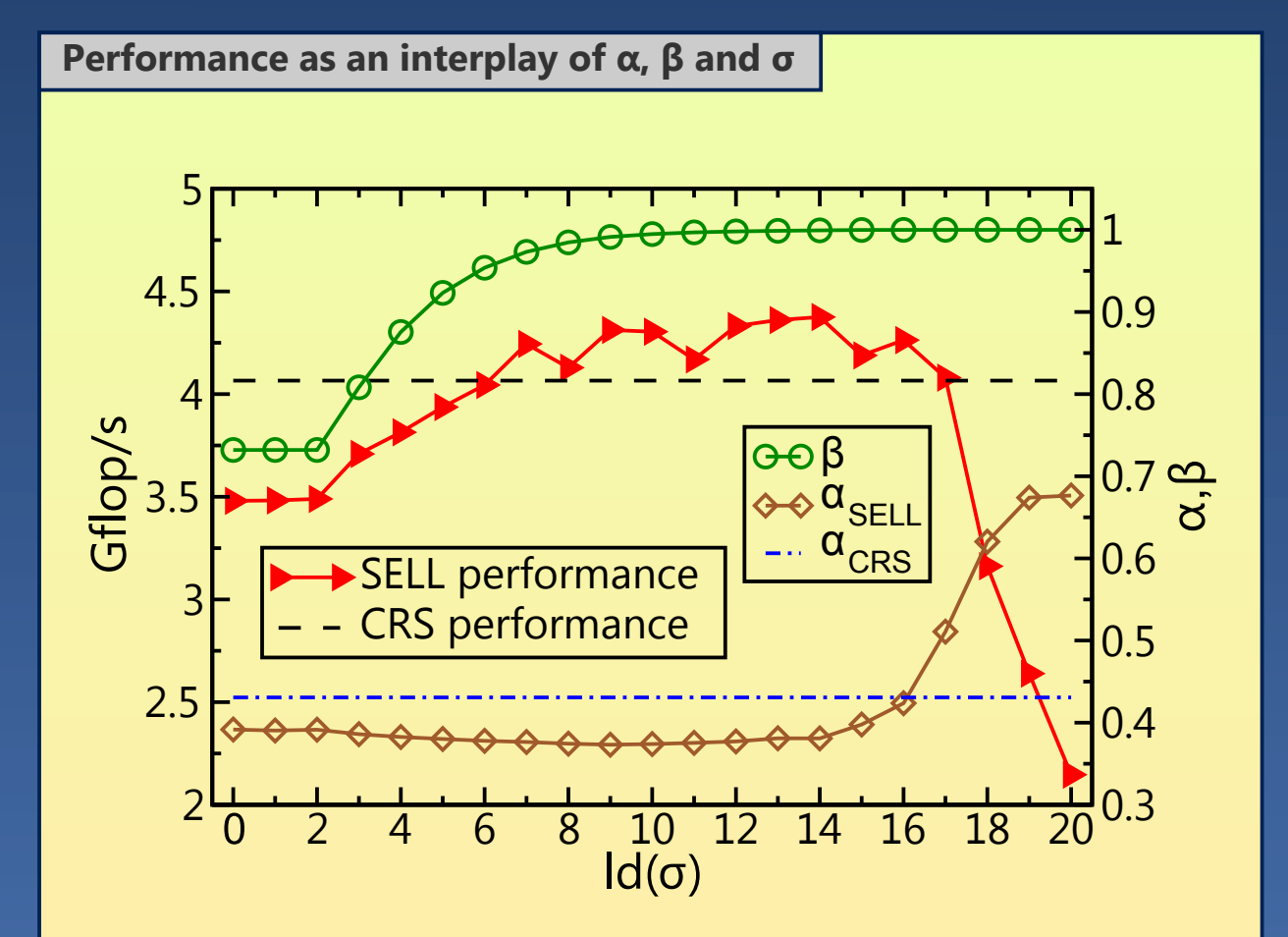
$$B_{\text{SELL}}^{\text{DP}} = \left( \frac{1}{\beta} \left( \frac{8+4}{2} \right) + \frac{x}{8\alpha+16} \frac{y}{N_{nz}} \right) \frac{\text{bytes}}{\text{flop}}$$

This knowledge can, e.g., be used for performance predictions based on a roofline model.

The factor  $\alpha$  quantifies the efficiency of the access to  $x$ .

- Each load of  $x$  goes to main memory:  $\alpha=1$  (or even larger if the cache is organized in cache lines)
- Access to  $x$  is perfectly cached:  $\alpha=1/N_{nz}$ .

Increasing the sorting scope  $\sigma$  results in a perfect  $\beta=1$ . However, if  $\sigma$  gets too large there is a high chance that the access pattern to  $x$  worsens ( $\alpha$  increases)  $\rightarrow$  performance drops.

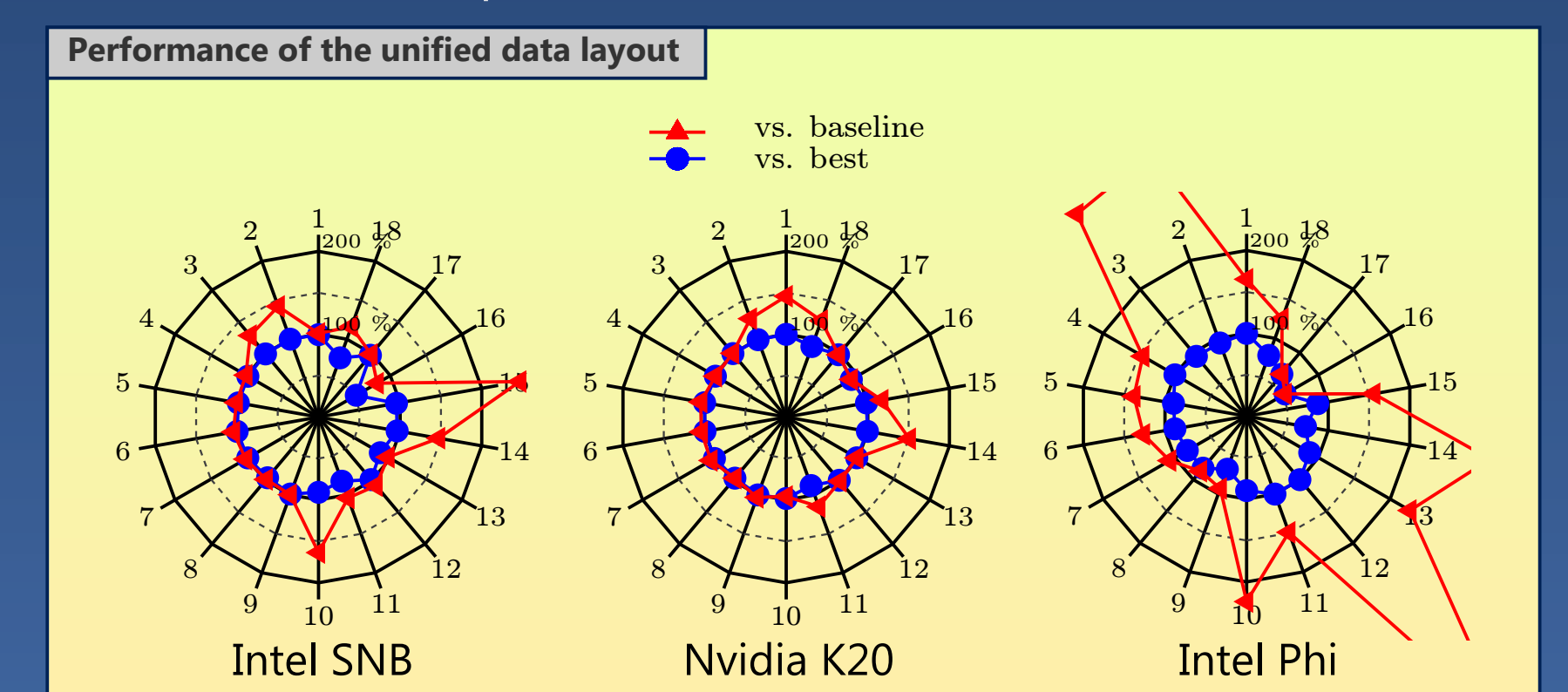


## Unified Layout Performance

When using SELL-C- $\sigma$  as a unified layout for the three architectures, we set  $C=32$  (warp size on Nvidia Kepler).  $\sigma$  is matrix-dependent.

The following image shows the SELL-32- $\sigma$  performance for 18 test cases relative to the

- baseline (same data as in Sec. Contribution) and
- best measured performance (best format,  $C$ , and  $\sigma$ )



The performance of the unified data layout barely falls behind the best performance significantly. For some cases with poor performance, better tuning parameters or a change to SELL-C-T- $\sigma$  (see Sec. Outlook) may improve the results.

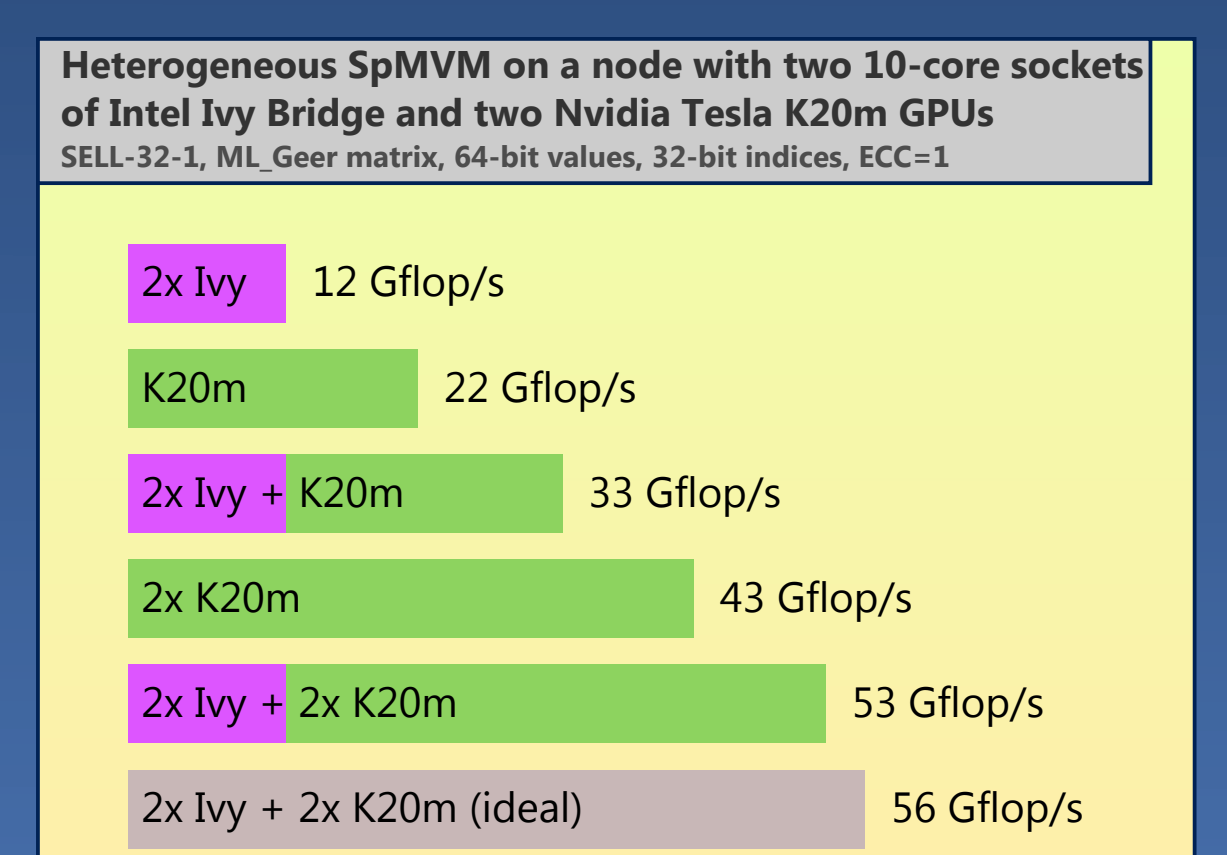
## Conclusions

We have presented SELL-C- $\sigma$  as an ideal candidate for a unified sparse matrix data layout. Compared to device-specific formats for three different architectures, it yields comparable or better performance on a wide range of sparse matrices. Especially on the Xeon Phi it is far superior to the standard CRS format.

## Outlook

The following issues are subject to further investigation:

- Increase the SELL-C- $\sigma$  performance for the critical test cases  $\rightarrow$  implement SELL-C-T- $\sigma$  where  $T$  threads run in a row
- Dynamically load-balanced heterogeneous spMVM  $\rightarrow$  Transfer chunks on demand to the accelerator  $\rightarrow$  Investigate offload-models like OpenACC  $\rightarrow$  Statically load-balanced implementation: see image
- Investigate SELL-C- $\sigma$  suitability for other numerical kernels besides SpMVM
- Find optimal  $\sigma$  automatically when reading in the matrix
- Exploit special matrix structures



## Acknowledgments

This work was supported (in part) by the German Research Foundation (DFG) through the Priority Programme 1648 "Software for Exascale Computing" (SPPEXA) under project ESSEX.



M. Kreutzer, G. Hager, G. Wellein, H. Fahske, and A. Bishop. A unified sparse matrix data format for modern processors with wide SIMD units. CoRR abs/1307.6209

D. R. Kincaid, T. C. Oppe and D. M. Young. ITPACKV 2D user's guide. Report CNA-232. The University of Texas at Austin, May 1989

N. Bell and M. Garland. Implementing sparse matrix-vector multiplication on throughput-oriented processors. In: Proceedings SC '09. DOI 10.1145/1654059.1654078

A. Monakov, A. Lokhtov, and A. Avetisyan. Automatically Tuning Sparse Matrix-Vector Multiplication for GPU Architectures. In: Y. Patt et al. (eds.), vol. 5952 of LNCS. DOI 10.1007/978-3-642-11515-8\_10

E. Saule, K. Kaya and Ü. V. Catalyurek. Performance evaluation of sparse matrix multiplication kernels on Intel Xeon Phi. CoRR abs/1302.1078

