

# Performance Engineering of the Kernel Polynomial Method on Large-Scale CPU-GPU Systems

Moritz Kreutzer<sup>\*</sup>, Georg Hager<sup>\*</sup>, Gerhard Wellein<sup>\*</sup>,  
Andreas Pieper<sup>#</sup>, Andreas Alvermann<sup>#</sup>, Holger Fehske<sup>#</sup>

<sup>\*</sup>: Erlangen Regional Computing Center, Friedrich-Alexander University of Erlangen-Nuremberg

<sup>#</sup>: Institute of Physics, Ernst Moritz Arndt University of Greifswald

29th IEEE International Parallel & Distributed Processing Symposium (IPDPS)  
May 26, 2015  
Hyderabad, India

This work was supported (in part) by the German Research Foundation (DFG) through the Priority Programs 1648 “Software for Exascale Computing” under project ESSEX and 1459 “Graphene”.

# Prologue

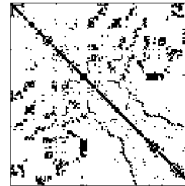
What is this about?



# The Kernel Polynomial Method (KPM)

Approximate the complete eigenvalue spectrum of a large sparse matrix.

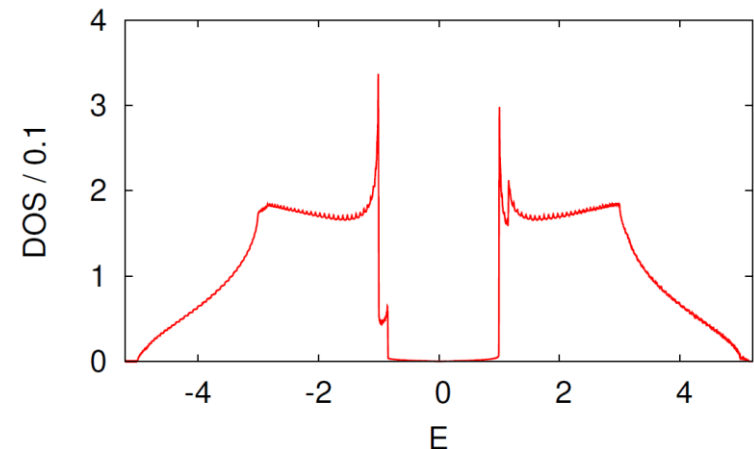
Large,  
Sparse



$$\mathbf{H} \mathbf{x} = \lambda \mathbf{x}$$

$$\{\lambda_1, \lambda_2, \dots, \lambda_k, \dots, \lambda_{n-1}, \lambda_n\}$$

Good approximation to full spectrum  
(e.g. Density of States)

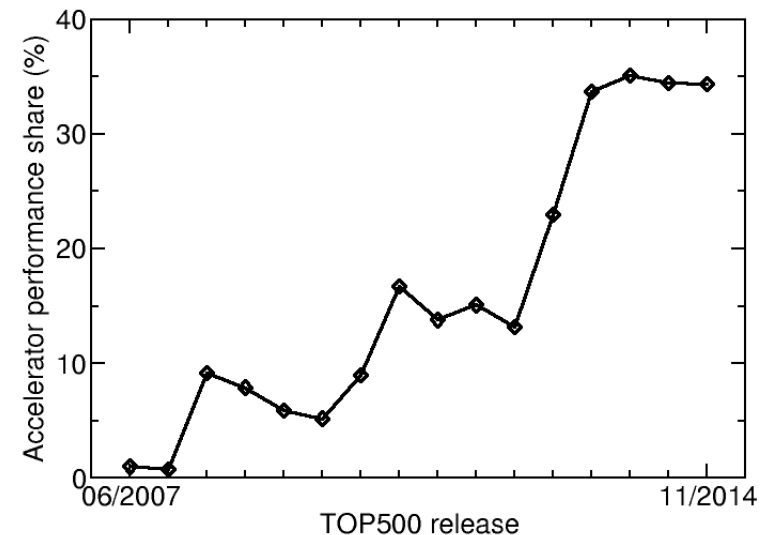


# Why optimize for heterogeneous systems?

One third of TOP500 performance stems from accelerators.

But: Few truly heterogeneous software.

(Using both CPUs and accelerators.)



# The Kernel Polynomial Method

## Algorithmic Analysis



# The Kernel Polynomial Method

Compute Chebyshev polynomials and moments.

## Basic algorithm and algorithmic optimizations: Exploit knowledge from all software layers!

**for**  $r = 0$  to  $R - 1$  **do**      Application: Loop over random initial states

$|v\rangle \leftarrow |\text{rand}()\rangle$

Initialization steps and computation of  $\eta_0, \eta_1$

**for**  $m = 1$  to  $M/2$  **do**      Algorithm: Loop over moments

$\text{swap}(|w\rangle, |v\rangle)$

$|u\rangle \leftarrow H|v\rangle$

$|u\rangle \leftarrow |u\rangle - b|v\rangle$

$|w\rangle \leftarrow -|w\rangle$

$|w\rangle \leftarrow |w\rangle + 2a|u\rangle$

$\eta_{2m} \leftarrow \langle v|v\rangle$

$\eta_{2m+1} \leftarrow \langle w|v\rangle$

**end for**

**end for**

Building blocks:  
(Sparse) linear  
algebra library

$-1$  **do**      steps and computation of  $\eta_0, \eta_1$

$M/2$  **do**

$\triangleright \text{spmv}(\cdot, |v\rangle)$

$\triangleright \text{axpy}(\cdot, (H - bI)|v\rangle - |w\rangle \ \&$

$\triangleright \text{scal}(\cdot) = \langle v|v\rangle \ \&$

$\triangleright \text{axpy}(\cdot) = \langle w|v\rangle$

$\triangleright \text{aug\_spmv}(\cdot)$

$\triangleright \text{nrm2}(\cdot) = \langle w|v\rangle$

$\triangleright \text{dot}(\cdot)$

Augmented Sparse

Matrix Vector Multiply

Dot Product

Multiple Vector Multiply

# Analysis of the Algorithmic Optimization

- Minimum code balance of vanilla algorithm:**

complex double precision values, 32-bit indices, 13 non-zeros per row, application: topological insulators

$$B_{vanilla} = 3.39 \text{ Bytes/Flop} \quad (B = \text{inverse computational intensity})$$

- Identified bottleneck: Memory bandwidth**

➔ **Decrease memory transfers to alleviate bottleneck**

- Algorithmic optimizations reduce code balance:**

$$B_{aug\_spmv} = 2.23 \text{ B/F} \quad \text{kernel fusion}$$

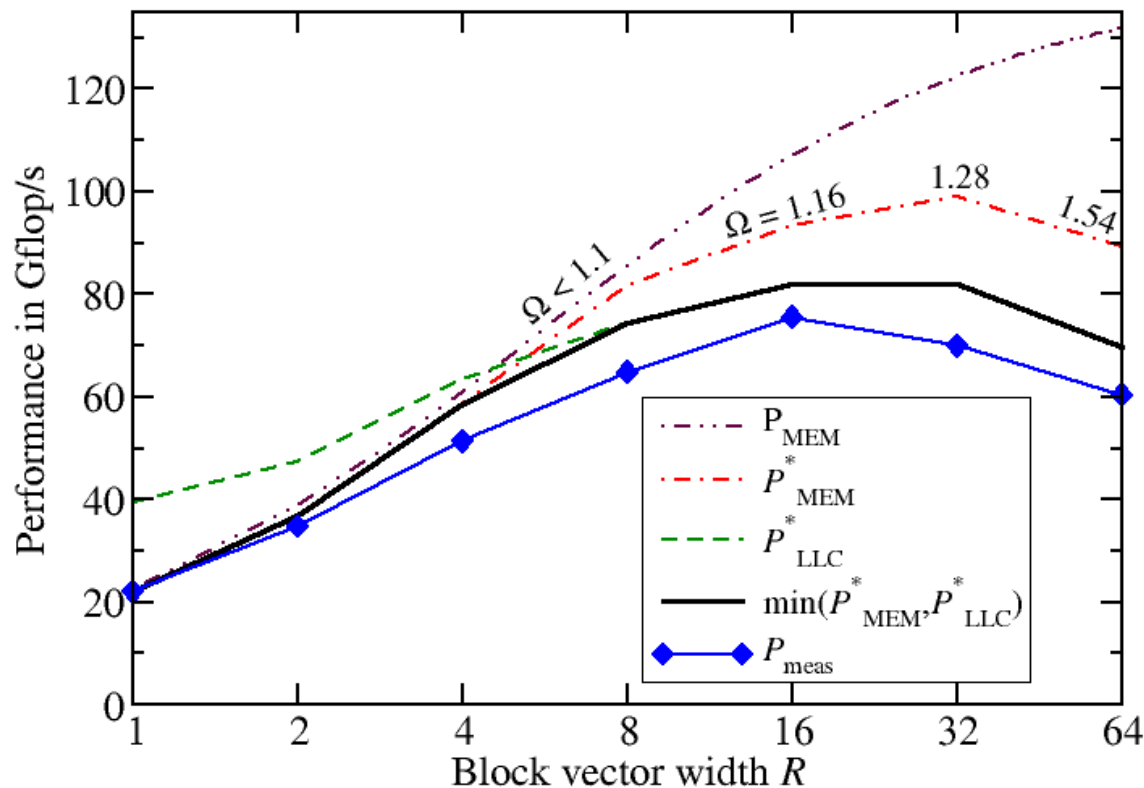
$$B_{aug\_spmmv}(R) = 1.88/R + 0.35 \text{ B/F} \quad \text{put } R \text{ vectors in block}$$

# Consequences of Algorithmic Optimization

- Mitigation of the relevant bottleneck  
→ Expected speedup 😊
- Other bottlenecks become relevant  
→ Achieved speedup may not be  $B_{vanilla}/B_{aug\_spmmv}$  😐
- Block vectors are best stored interleaved  
→ May impose larger changes to the codebase 😐
- `aug_spmmv()` no part of standard libraries  
→ Implementation by hand is necessary 😐



# CPU roofline performance model



$$P = \frac{b}{B} \text{ Gflop/s}$$

→ Performance limit for bandwidth-bound code

$b$  = max. bandwidth = 50 GB/s  
 $B$  = code balance

$$\Omega = \frac{\text{Actual data transfers}}{\text{Minimum data transfers}}$$

S. Williams, A. Waterman, D. Patterson: “Roofline: An insightful visual performance model for multicore architectures”, *Commun. ACM*, vol. 52, p. 65, 2009.

# Implementation

How to harness a heterogeneous machine in an efficient way?



# Implementation

Algorithmic optimizations lead to a *potential* speedup.

→ We “merely” need an efficient implementation!

## Data or task parallelism?

- **MAGMA: task parallelism between devices**  
<http://icl.cs.utk.edu/magma/>
- **Kernel fusion**  **Task parallelism**

→ Data-parallel approach suits our needs

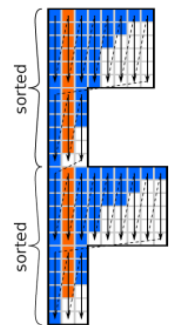
# Implementation

## Data-parallel heterogeneous work distribution

- Static work-distribution by matrix rows/entries
- Device workload  $\leftrightarrow$  device performance

## SELL-C- $\sigma$ sparse matrix storage format

- Unified format for all relevant devices
- Currently no runtime-exchange of matrix data (dynamic load balancing, future work)



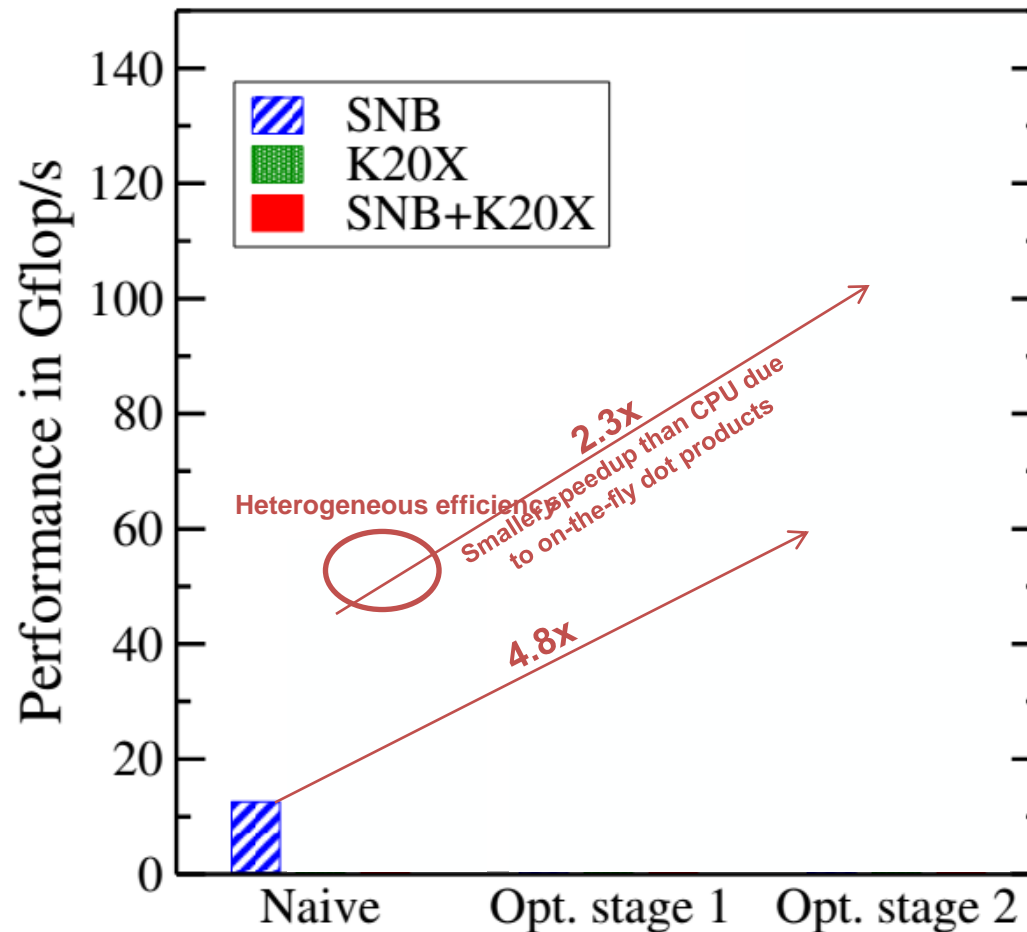
M. Kreutzer, G. Hager, G. Wellein, H. Fehske, A. R. Bishop, “A unified sparse matrix data format for efficient general sparse matrix-vector multiplication on modern processors with wide SIMD units”, *SIAM J. Sci. Comput.*, vol. 36, p. C401, 2014

# Performance results

Does all this really pay off?

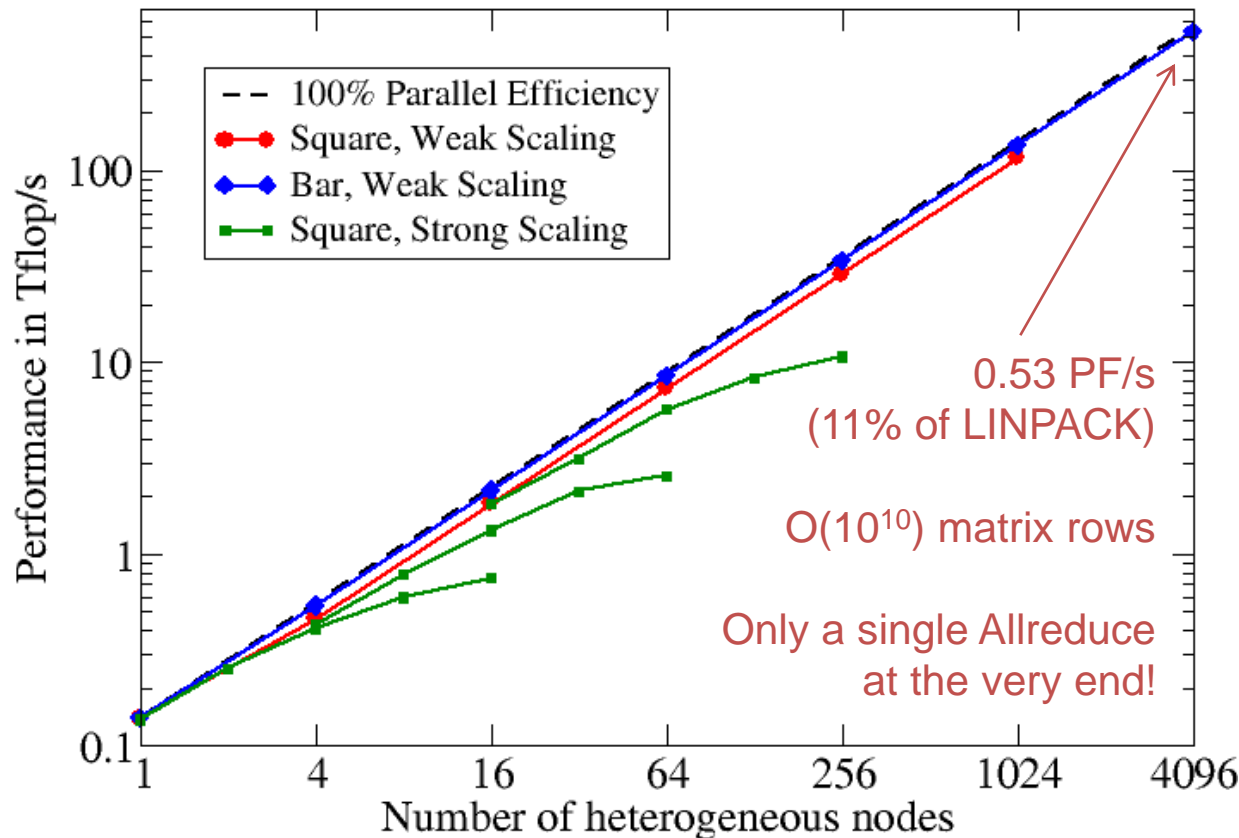


# Single-node Heterogeneous Performance



SNB: Intel Xeon Sandy Bridge, K20X: Nvidia Tesla K20X, Complex double precision matrix/vectors (topological insulator)

# Large-scale Heterogeneous Performance



## CRAY XC30 – Piz Daint\*



- 5272 nodes, each w/
  - 1 octacore Intel Sandy Bridge
  - 1 Nvidia Kepler K20x
- Peak: 7.8 Pflop/s
- LINPACK: 6.3 Pflop/s
- Largest system in Europe

\*Thanks to CSCS/O. Schenk/T. Schulthess for granting access and compute time

# Epilogue

Try it out! (If you want...)





# Download our building block library and KPM application: <http://tiny.cc/ghost>



*General, Hybrid, and Optimized Sparse Toolkit*

- **MPI + OpenMP + SIMD + CUDA**
- **Transparent data-parallel heterogeneous execution**
- **Affinity-aware task parallelism (checkpointing, comm. hiding, etc.)**
- **Support for block vectors**
  - Automatic code generation for common block vector sizes
  - Hand-implemented tall skinny dense matrix kernels
- **Fused kernels (arbitrarily “augmented SpMMV”)**
- **SELL-C- $\sigma$  heterogeneous sparse matrix format**
- **Various sparse eigensolvers implemented and downloadable**
- . . .

# Backup Slides

Only in the unlikely case I was too fast...



# Conclusions

- **Model-guided performance engineering of KPM on CPU and GPU**
- **Decoupling from main memory bandwidth**
- **Optimized node-level performance**
- **Embedding into massively-parallel application code**
- **Fully heterogeneous peta-scale performance for a sparse solver**

# Outlook

- **Applications besides KPM**
- **Automatic (& dynamic) load balancing**
- **Optimized GPU-CPU-MPI communication**
- **Further optimization techniques (cache blocking, ...)**
- **Performance engineering for Xeon Phi (already supported)**