

Efficient Eigensolver Algorithms on Accelerator Based Architectures

Stan Tomov

w/ Azzam Haidar, Piotr Luszczek, Ichitaro Yamazaki, Mark Gates, and Jack Dongarra

Innovative Computing Laboratory
Department of Computer Science
University of Tennessee, Knoxville

SIAM LA 2015
Atlanta, GA
October 29, 2015

Outline

- **Motivation**
- **Eigenvalue and SVD problems**
- **Main challenges**
 - **Parallelism**
 - **Heterogeneity**
 - **Communication**
- **One and two-stage algorithms**
- **Conclusions and future directions**

Linear algebra in applications

Dense Linear Algebra (DLA) is needed in a wide variety of science and engineering applications:

- **Linear systems:** **Solve $Ax = b$**
 - Computational electromagnetics, material science, applications using boundary integral equations, airflow past wings, fluid flow around ship and other offshore constructions, and many more
- **Least squares:** **Find x to minimize $\|Ax - b\|$**
 - Computational statistics (e.g., linear least squares or ordinary least squares), econometrics, control theory, signal processing, curve fitting, and many more
- **Eigenproblems:** **Solve $Ax = \lambda x$**
 - Computational chemistry, quantum mechanics, material science, face recognition, PCA, data-mining, marketing, Google Page Rank, spectral clustering, vibrational analysis, compression, and many more
- **Singular Value Decomposition (SVD): $A = U \Sigma V^*$ ($Au = \sigma v$ and $A^*v = \sigma u$)**
 - Information retrieval, web search, signal processing, big data analytics, low rank matrix approximation, total least squares minimization, pseudo-inverse, and many more
- **Many variations depending on structure of A**
 - A can be symmetric, positive definite, tridiagonal, Hessenberg, banded, sparse with dense blocks, etc.
- **DLA is crucial to the development of sparse solvers**

General Overview: the Eigen-problem algorithms

1. Singular Value Decomposition $A = U\Sigma V^T$
 - Bi-Diagonalization Reduction + solve + back transformation.
2. Symmetric Eigenvalue Problem $Ax = \lambda x$
 - Tri-Diagonalization Reduction + solve + back transformation.
3. Non-symmetric Eigenvalue Problem $Ax = \lambda x$
 - Hessenburg Reduction + solve + back transformation.

General Overview: the Eigen-problem algorithms

■ Three phases:

- Reduction phase – reduce A to condensed form



$$A = Q H Q^*$$

Nonsymmetric
Eigenvalue problem



$$A = Q T Q^*$$

Symmetric
Eigenvalue problem



$$A = Q B P^*$$

Singular value
Decomposition
problem

- Solution phase – compute the eigenvalues, eigenvectors and singular value of the condensed matrix

$$H = E T E^*$$

$$\implies A = Q (E T E^*) Q^*$$

$$T = E \Lambda E^*$$

$$\implies A = Q (E \Lambda E^*) Q^*$$

$$B = X S Y^*$$

$$\implies A = Q (X S Y^*) P^*$$

- Back transformation phase – Update the computed eigenvectors and orthogonal matrix

$$Z = Q E$$

$$Z = Q E$$

$$U = Q X$$

$$V^* = Y^* P^*$$

Major challenges to software on heterogeneous architectures

- **Parallelism**





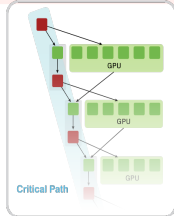
- **Explosion in parallelism; stalled clock rates**
[2x1,664 CUDA cores in K80]

- **Heterogeneity**

- **A few heavyweight and many lightweight cores**
- **Multiple processor architectures, memory systems, etc.**

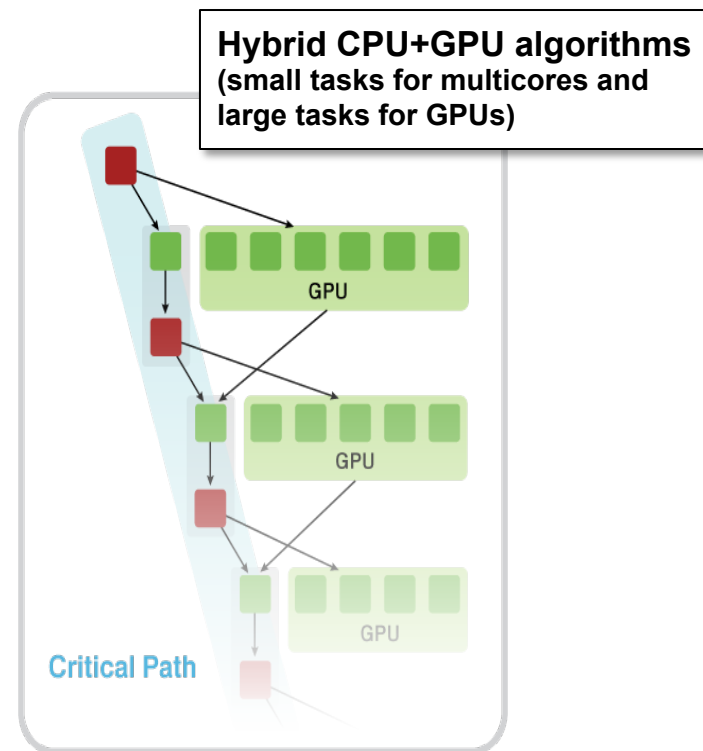
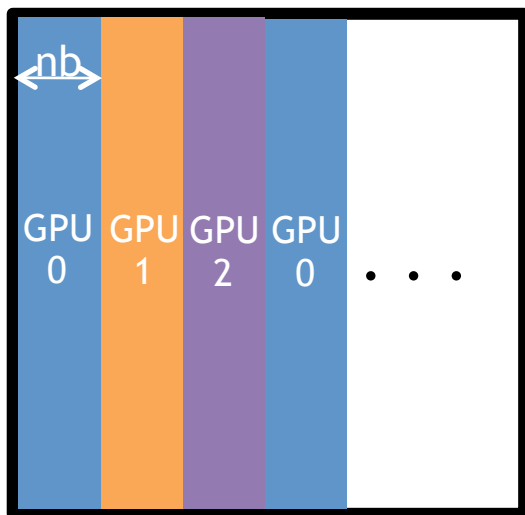
- **Communications**

- **Growing gap between peak performance and bandwidth**

Software/Algorithms follow hardware evolution in time		
LINPACK (70's) (Vector operations)		Rely on - Level-1 BLAS operations
LAPACK (80's) (Blocking, cache friendly)		Rely on - Level-3 BLAS operations
ScaLAPACK (90's) (Distributed Memory)		Rely on - PBLAS Mess Passing
PLASMA (00's) New Algorithms (many-core friendly)		Rely on - a DAG/scheduler - block data layout - some extra kernels
MAGMA Hybrid Algorithms (heterogeneity friendly)		Rely on - hybrid scheduler - hybrid kernels (for nested parallelism)

Parallelism

- Algorithms expressed in terms of BLAS
- Task based approach with event/data driven execution
 - Algorithms as DAGs
 - Runtime system schedules the execution (static, QUARK, PaRSEC, OpenMP4, SMPs, StarPU, etc.)
- **1-D block-cyclic distribution on single node** (or 2-D block-cyclic on distributed memory systems)



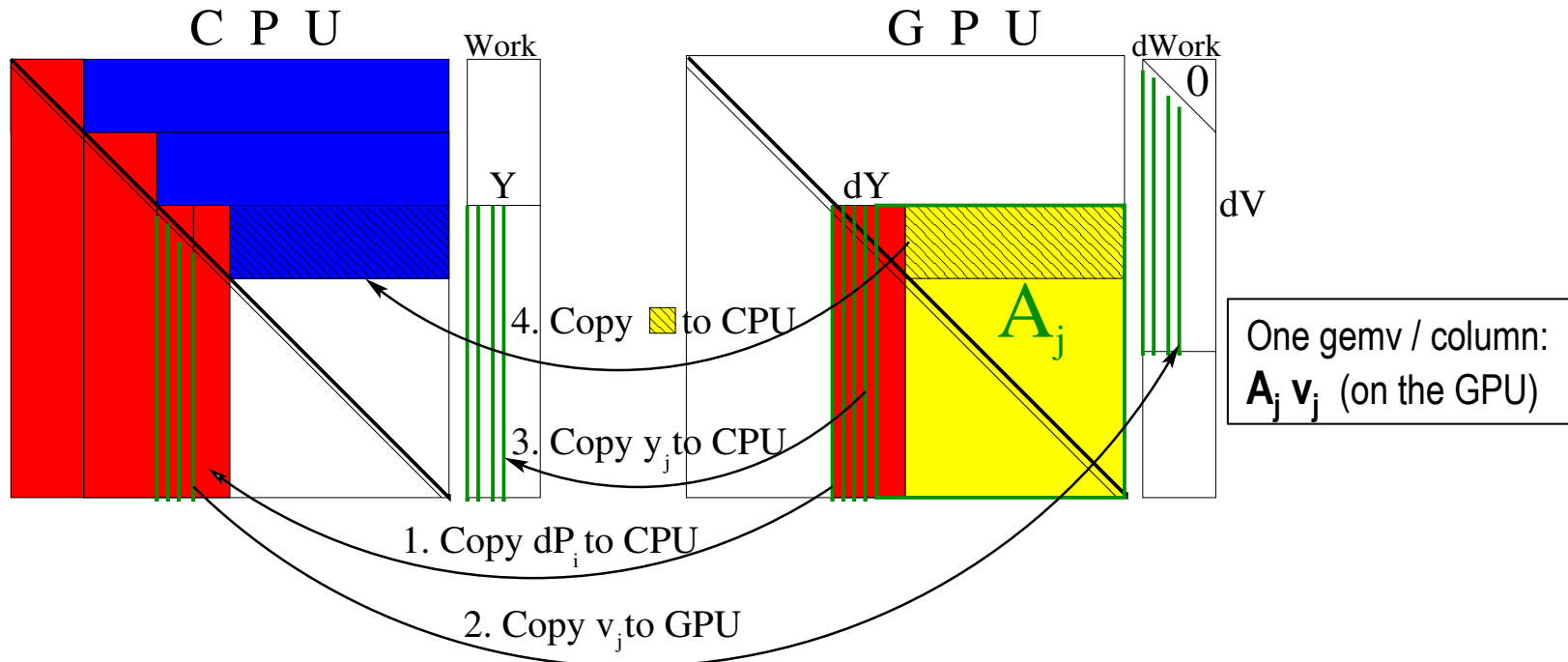
Heterogeneity ...

- **Example:**

- Panels are in general scheduled on the CPU
- For the eigen-problems and SVD panels can be hybrid

Hybrid panel in the Hessenberg reduction

(Next householder vector from the panel factorization is sent to GPU and multiplied by trailing matrix)

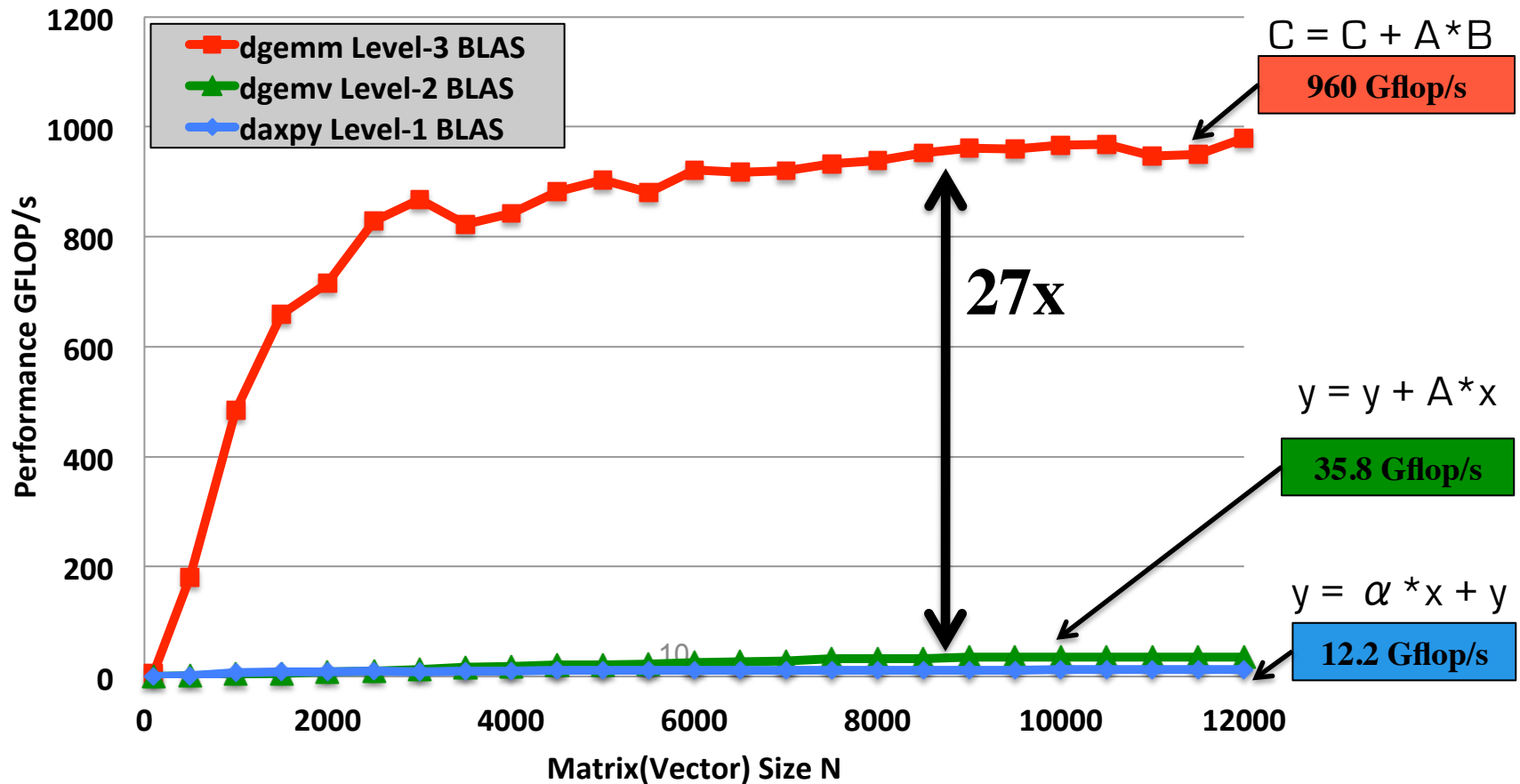


Communication

Not all flops are equal (some need more data movements)

Level 1, 2 and 3 BLAS

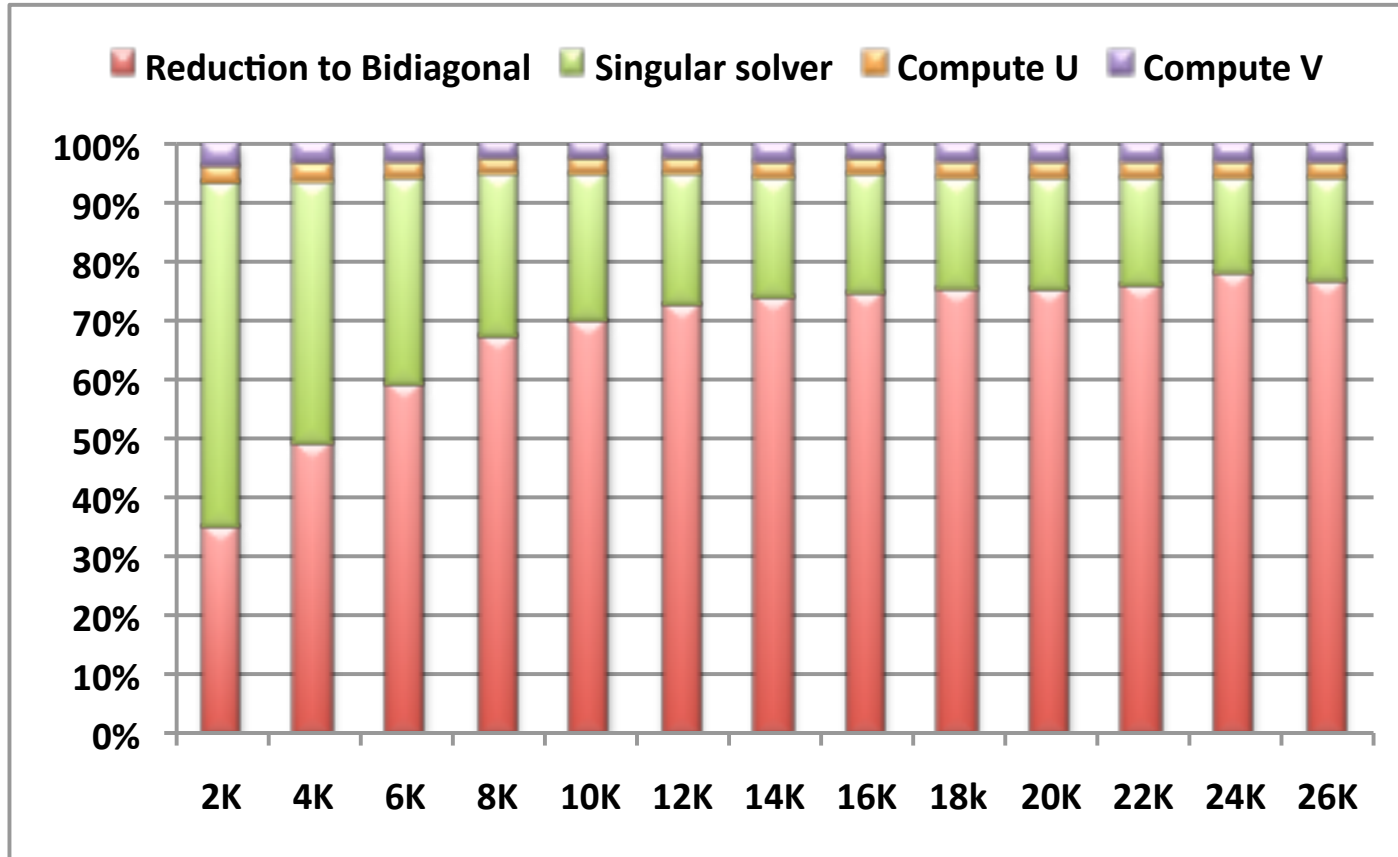
60 core Intel Xeon Phi KNC 7120, 1.23 GHz, Peak DP = 1208 Gflop/s



Profiling DGESDD (CPU) Classical Reduction and D&C

Three Phases:

1. Reduction
2. Solver
3. Back transform



}	Back-trans	$4n^3$
}	Solver	$8n^3$
}		$\frac{8n^3}{3}$
}	Reduction	$8n^3$
}		$\frac{8n^3}{3}$
}	Total flops	$\approx 9n^3$

- DGESDD – computes singular value decomposition of $m \times n$ matrix
- Reduction to Bidiagonal form(DGEBRD) – $\approx 75\%$ of the total computation cost

Communication ...

Bidiagonal reduction

Performance Bound Analysis

- Total cost for bidiagonal reduction (1st phase of SVD)

$$\approx 4n_b \sum_{n_b}^{n/n_b} l^2 + 4n_b \sum_{2n_b}^{n-n_b} k^2$$

$$\approx \frac{8}{3}n^3$$

$$\approx \frac{4}{3}n^3_{gemv} + \frac{4}{3}n^3_{gemm}$$

- Performance upper bound

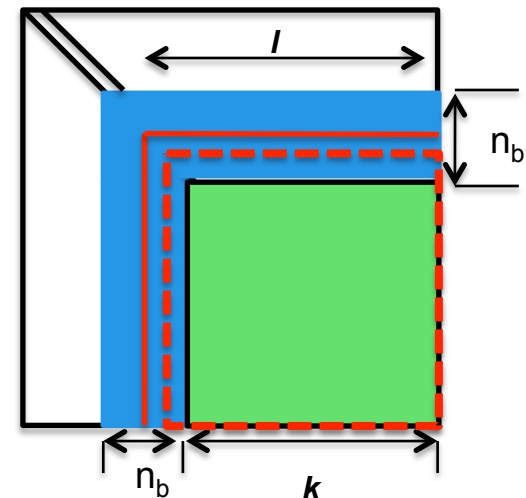
$$P_{max} = \frac{\text{number of operations}}{\text{minimum time } t_{min}}$$

$$= \frac{\frac{8}{3}n^3}{\frac{4}{3}n^3 * \frac{1}{P_{gemv}} + \frac{4}{3}n^3 * \frac{1}{P_{gemm}}}$$

$$= \frac{2 * P_{gemm} * P_{gemv}}{P_{gemm} + P_{gemv}}$$

$$\approx 2P_{gemv}$$

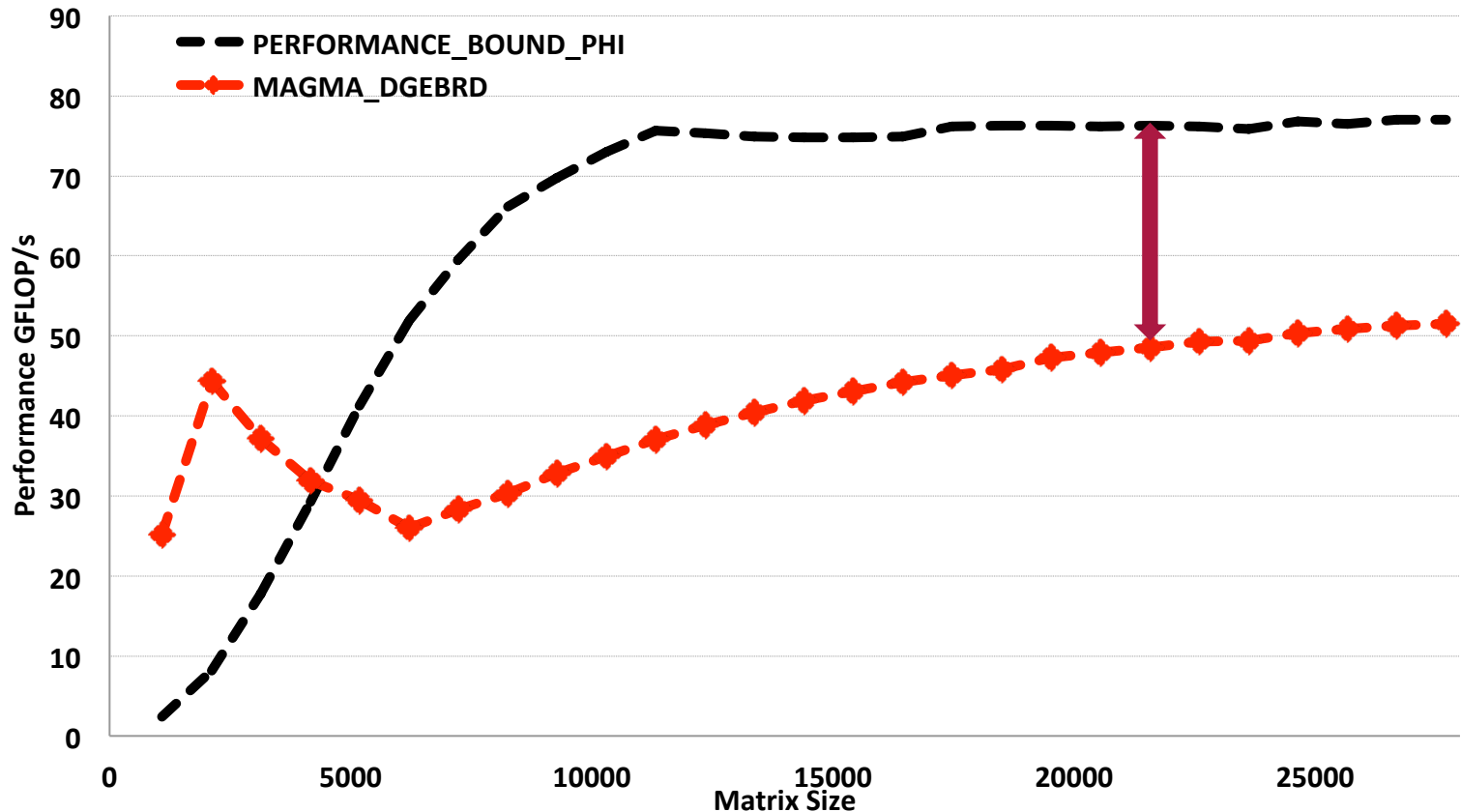
$$\text{when } P_{gemm} \gg P_{gemv}.$$



- Trailing Matrix Update – Level 3 BLAS update
- Panel - Level 2 BLAS update **GEMV**
- Panel

Communication ...

Bidiagonal reduction – Performance (Xeon Phi)



- **Bottleneck:**

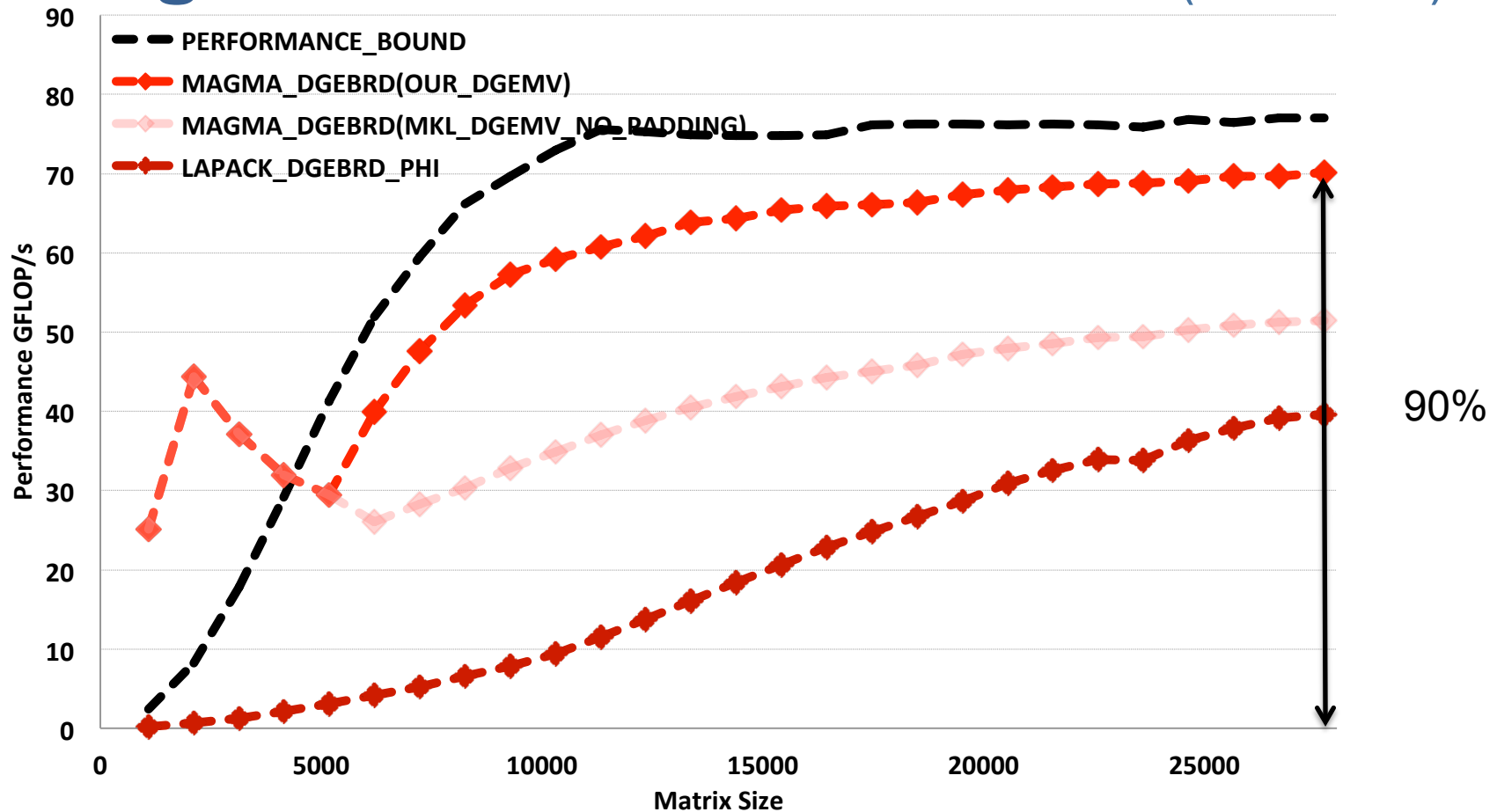
- Accelerator implementation reaches **60% of the performance bound** (Theoretical peak > 1 Tflop/s)
- The gap indicates that DGEMV inside DGEBRD does not perform well.

- **Observations:**

MKL's DGEMV performance depends on both, the size and the memory alignment of the matrix

Communication ...

Bidiagonal reduction – Performance (Xeon Phi)



- Our implementation is within 90% of the performance bound
- Hierarchical communications (blocking + prefetching for both L1 and L2 cache) is essential here

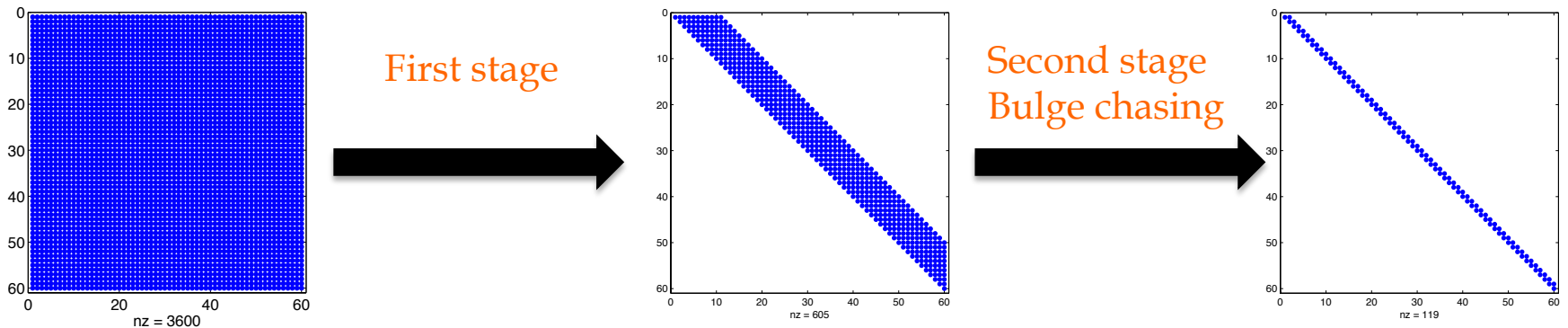
Communication ...

Bidiagonal reduction

- The only way to do better – design a new algorithm to communicate less

Communication ... Bidiagonal reduction

- The only way to do better – design a new algorithm to communicate less
- Recent work on developing new “2-stage” algorithms



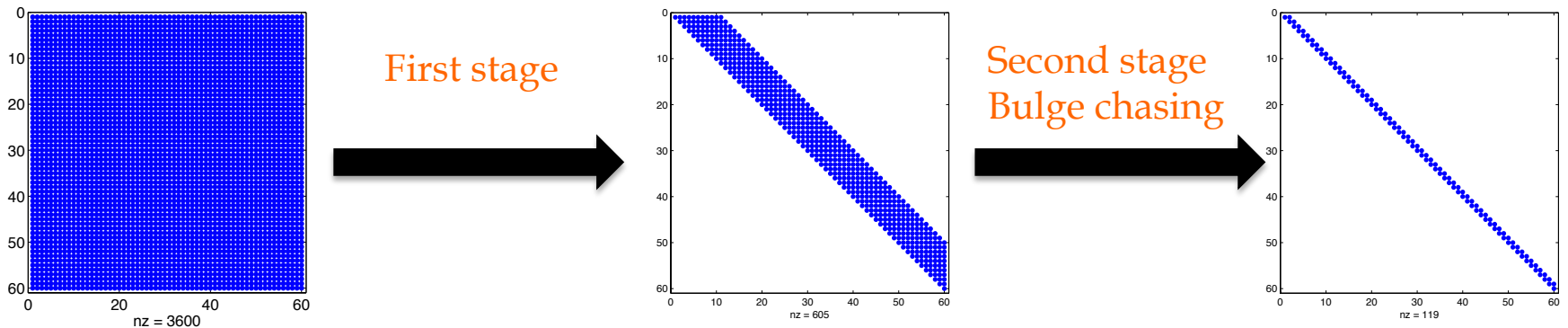
$$\begin{aligned}
 \text{flops} &\approx \sum_{s=1}^{n-n_b} 2n_b^3 + (nt-s)3n_b^3 + (nt-s)\frac{10}{3}n_b^3 + (nt-s) \times (nt-s)5n_b^3 \\
 &+ \sum_{s=1}^{n_b} 2n_b^3 + (nt-s-1)3n_b^3 + (nt-s-1)\frac{10}{3}n_b^3 + (nt-s) \times (nt-s-1)5n_b^3 \\
 &\approx \frac{10}{3}n^3 + \frac{10n_b}{3}n^2 + \frac{2n_b}{3}n^3 \\
 &\approx \frac{10}{3}n^3(\text{gemm})_{\text{first stage}}
 \end{aligned}$$

$$\text{flops} = 6 \times n_b \times n^2(\text{gemv})_{\text{second stage}}$$

More Flops, original did $\frac{8}{3} n^3$

Communication ... Bidiagonal reduction

- The only way to do better – design a new algorithm to communicate less
- Recent work on developing new “2-stage” algorithms



$$\text{speedup} = \frac{\text{time of one-stage}}{\text{time of two-stage}}$$

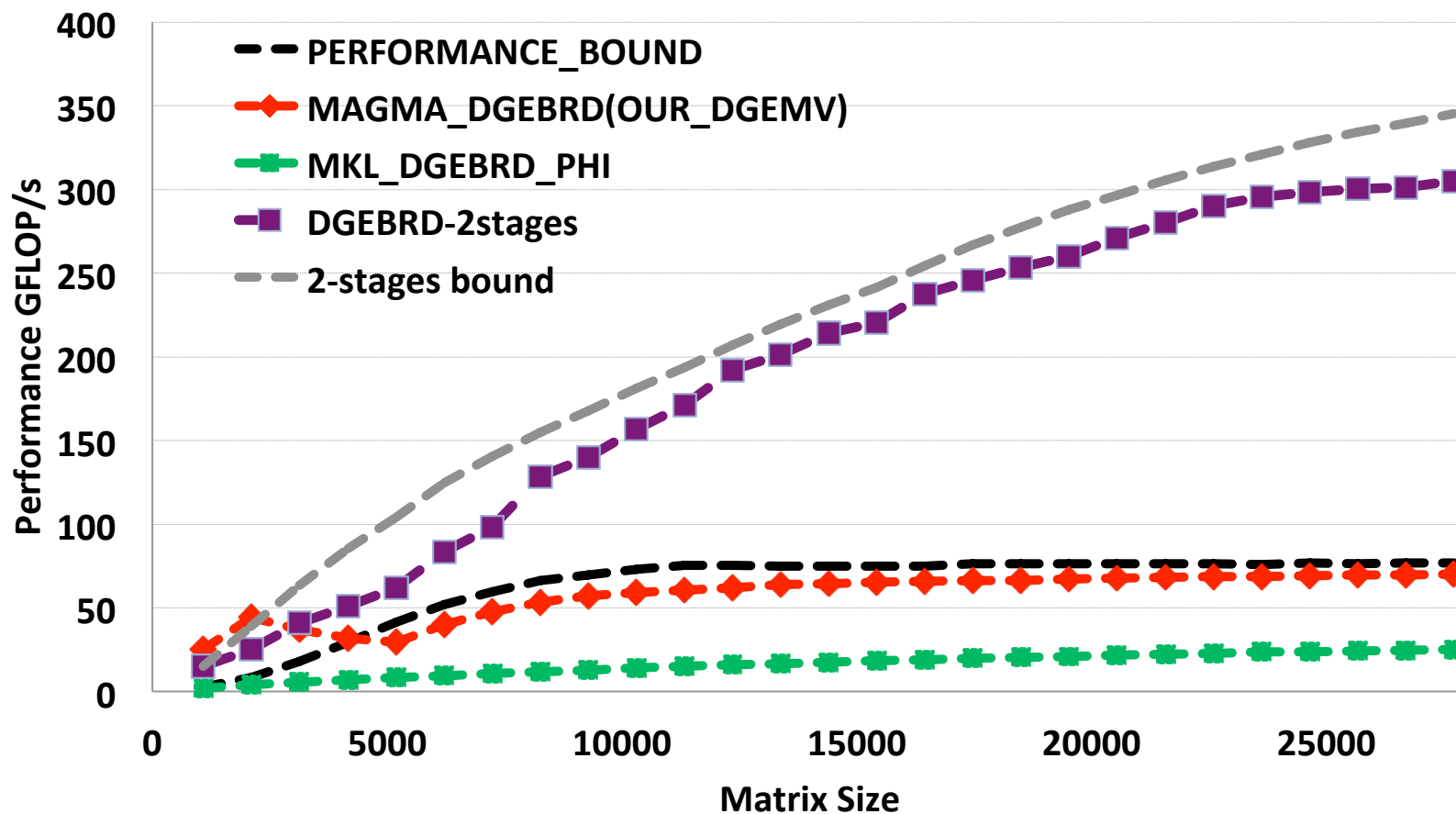
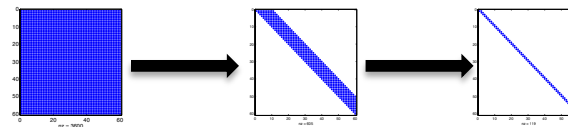
$$= \frac{4n^3/3P_{\text{gemv}} + 4n^3/3P_{\text{gemm}}}{10n^3/3P_{\text{gemm}} + 6n_b n^2/P_{\text{gemv}}}$$

$$\implies \frac{84}{70} \leq \text{Speedup} \leq \frac{84}{15}$$

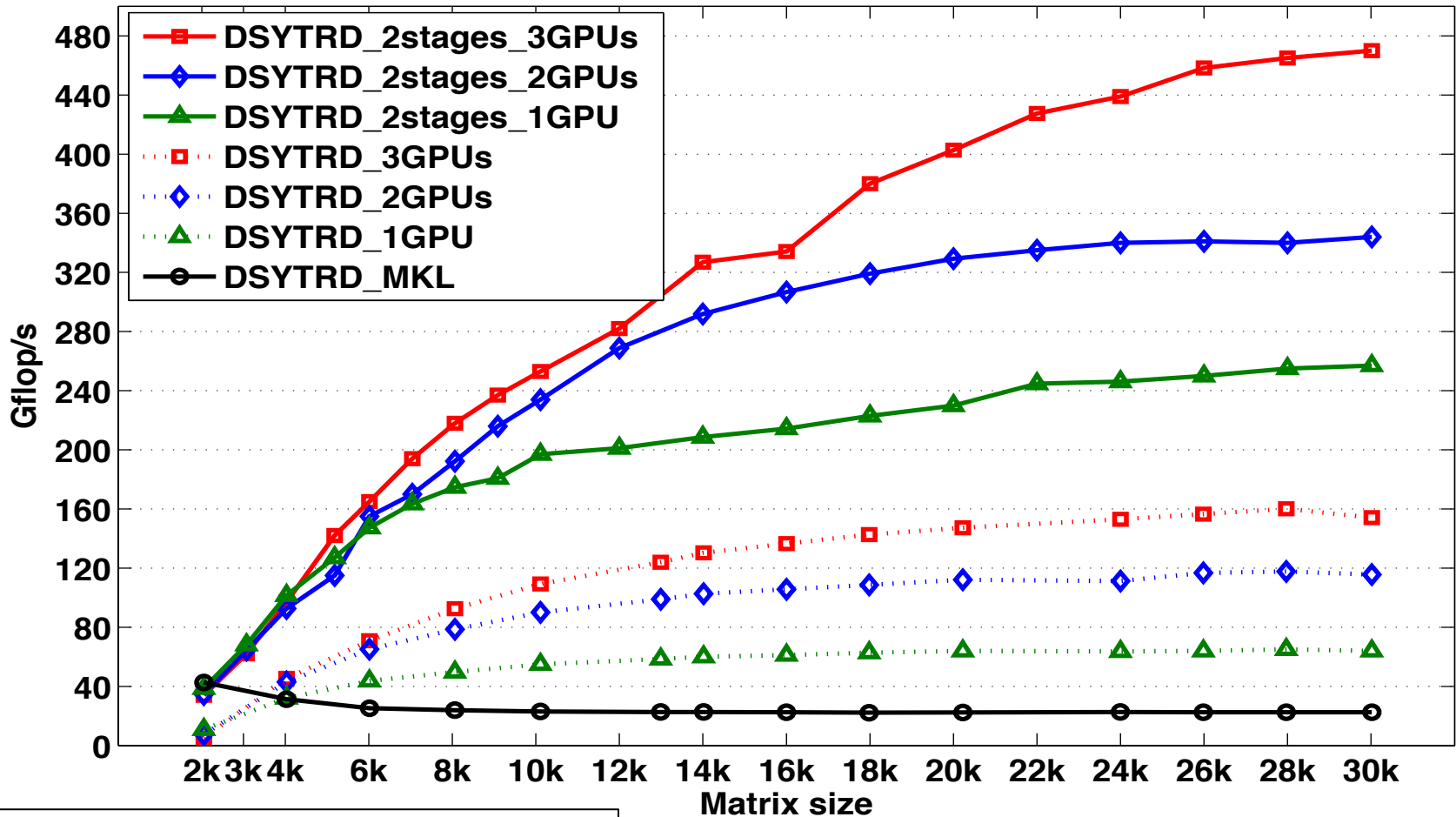
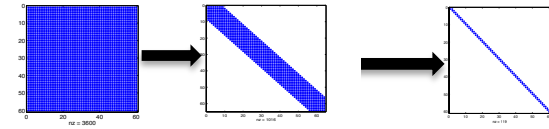
$$\implies 1.2 \leq \text{Speedup} \leq 5.6$$

if P_{gemm} is about 20x P_{gemv} and $120 \leq n_b \leq 240$.

Communication ... Bidiagonal reduction

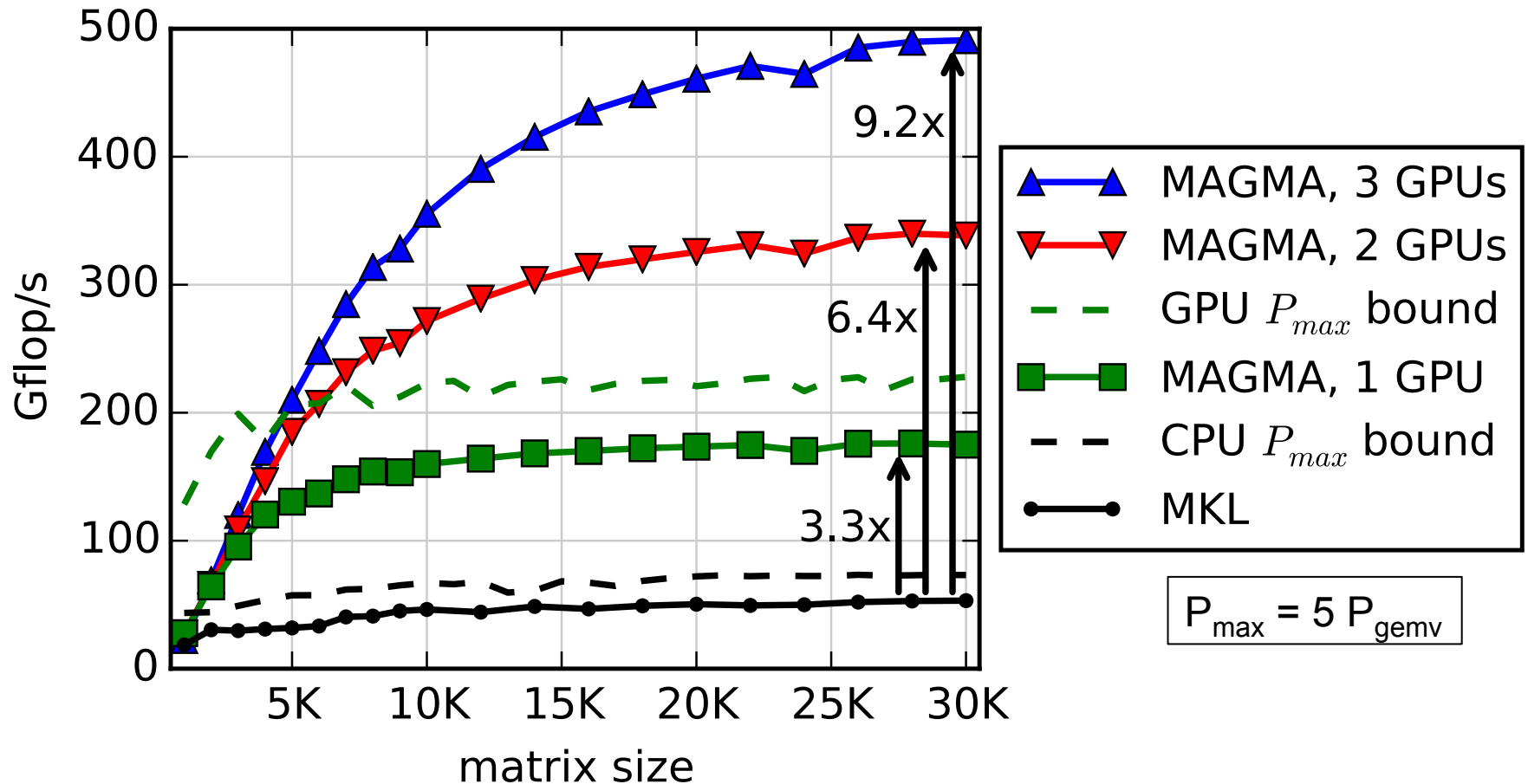


Communication ... Tridiagonal reduction



3 NVIDIA GPUs (M2090@ 1.1 GHz, 5.4 GB)
2 x 6 Intel Cores (X5660 @ 2.8 GHz, 23 GB)

Communication ... Hessenberg reduction (dgehrd)



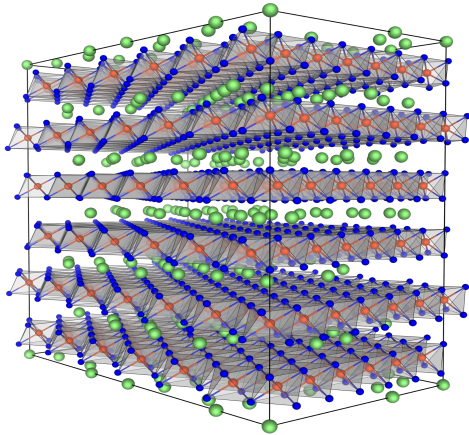
3 NVIDIA GPUs (K40@ 876 MHz)
2 x 8 Intel Cores (E5-2670 @ 2.6 GHz)

$$P_{max} = 5 P_{gemv}$$

Distributed memory systems

Symmetric eigensolvers

with Thomas Schulthess et al. CSCS



	Setup H2O	Solve HC=40C	The rest	total
28x28(2R:4T) ScaLAPACK	463.7	3839.7	61.6	4365.0
28x28(2R:4T) ELPA	471.7	1199.3	60.7	1731.7
14x14(1R:8T) MAGMA	166.7	911.9	79.9	1158.5

MAGMA is **1.5 times faster** than CPU

MAGMA is **2 times more energy efficient**

Cray XC30:

- 8-core Intel Xeon E5-2670 Sandy Bridge socket
- Nvidia K20X

- **R. Solca, A. Kozhevnikov, A. Haidar, S. Tomov, T. Schulthess, and J. Dongarra**, *Efficient implementation of quantum material simulations on distributed CPU-GPU systems*. SC'15, Best Paper Award finalist, Austin, TX, November 15-20, 2015.
- **A. Haidar, S. Tomov, J. Dongarra, T. Schulthess, and R. Solca**, *A novel hybrid CPU-GPU generalized eigensolver for electronic structure calculations based on fine grained memory aware tasks*, International Journal of High Performance Computing Applications , vol. 28, no 2, pp. 196—209, May 2014.
- **A. Haidar, R. Solca, S. Tomov, T. Schulthess and J. Dongarra**. *Leading edge multi-GPU algorithms for generalized eigenproblems for electronic structure calculations*. International Supercomputing Conference IEEE-ISC 2013.

Conclusions

- **Presented methodology and algorithms for efficient eigensolvers on accelerator based systems**
 - **Parallelism** (task based + runtime system)
 - **Heterogeneity** (unified support for various accelerators + scheduling)
 - **Communication** (both kernel and algorithmic level; one and two-stage algorithms)
- **Significant acceleration and energy efficiency compared to standard algorithms on CPUs**

Future Directions

- **GPU-only algorithms and implementations**
- **Further improvements regarding main challenges to meet emerging application demands**
- **Optimizations for small sizes, including batched versions**
- **Embedded systems**

Collaborators and Support

MAGMA team

<http://icl.cs.utk.edu/magma>



PLASMA team

<http://icl.cs.utk.edu/plasma>



Collaborating partners

University of Tennessee, Knoxville

University of California, Berkeley

University of Colorado, Denver

INRIA, France (StarPU team)

KAUST, Saudi Arabia



U.S. DEPARTMENT OF
ENERGY

