

Exploiting Batched Operation in Applications

David Keyes and Hatem Ltaief

Extreme Computing Research Center
King Abdullah University of Science and Technology

Workshop on Batched, Reproducible, and Reduced Precision BLAS
GATech

Feb 25, 2017



Outline

- 1 Motivations
- 2 Real Scientific Applications
- 3 Cholesky-based Matrix Computations
- 4 Climate/Weather Prediction Application
- 5 KBLAS

Outline

- 1 Motivations
- 2 Real Scientific Applications
- 3 Cholesky-based Matrix Computations
- 4 Climate/Weather Prediction Application
- 5 KBLAS

Motivations

- Covariance Matrix Problems
 - Ubiquitous in computational science and engineering
 - Symmetric, positive-definite matrix structure
 - (Apparently) Dense matrices
 - Often data-sparse
 - Decay of parameter correlations with distance
 - Hierarchically of low rank
 - **Convergence big data / HPC**
- Sparse direct and iterative solvers
 - Schur complement
 - Preconditioning

Outline

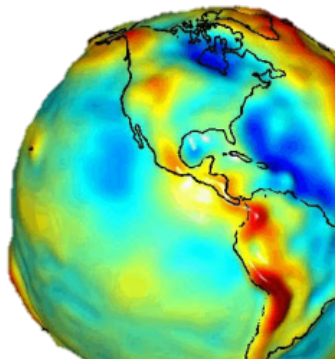
- 1 Motivations
- 2 Real Scientific Applications**
- 3 Cholesky-based Matrix Computations
- 4 Climate/Weather Prediction Application
- 5 KBLAS

Geospatial Statistics

- Multivariate large spatial data sets in climate/weather modeling to improve prediction

$$\ell(\boldsymbol{\theta}) = -\frac{1}{2}\mathbf{Z}^T \Sigma^{-1}(\boldsymbol{\theta})\mathbf{Z} - \frac{1}{2}\log|\Sigma(\boldsymbol{\theta})|$$

(a) Problem Definition.



(b) Temperature prediction.

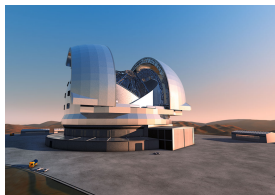
Figure: Climate/weather modeling.

Computational Ground-Based Astronomy

- Enhancing the observed image quality using MOAO by filtering out the noise coming from the adaptive optics instrumentation and the atmospheric turbulence.

$$R = C_{tm} \cdot C_{mm}^{-1} \quad C_{ee} = C_{tt} - C_{tm} R^t - R C_{tm}^t + R C_{mm} R^t$$

(a) Problem Definition.

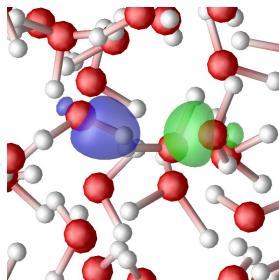


(b) The E-ELT.

Figure: Finding new galaxies.

Computing the Eigenspectrum for Symmetric Hierarchical Low Rank Matrix

- Structural and vibrational analysis to problems in computational physics and chemistry like electronic and band structure calculations



$$(A - \lambda B)x = 0$$

(a) Problem Definition. (b) Electronic structure.

Figure: Design of new materials.

Outline

- 1 Motivations
- 2 Real Scientific Applications
- 3 Cholesky-based Matrix Computations**
- 4 Climate/Weather Prediction Application
- 5 KBLAS

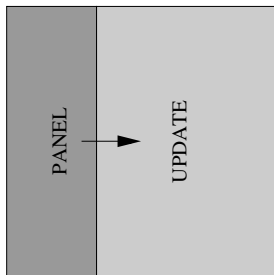
Matrix Form

The Cholesky factorization of an $N \times N$ real symmetric, positive-definite matrix A has the form

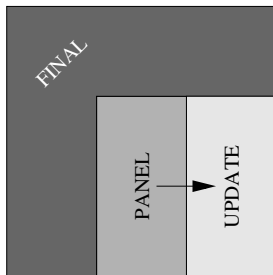
$$A = LL^T,$$

where L is an $N \times N$ real lower triangular matrix with positive diagonal elements.

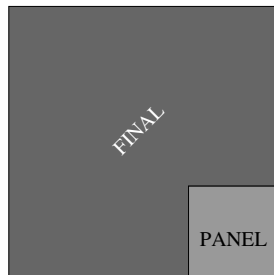
LAPACK Block Algorithms



(a) First step.



(b) Second step.



(c) Third step.

Figure: Block Algorithms.

PLASMA Tile Algorithms

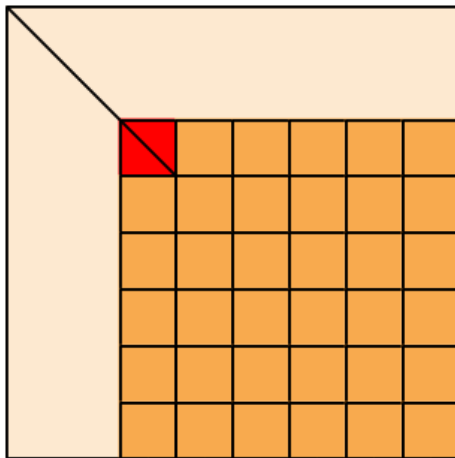


Figure: Tile Algorithms.

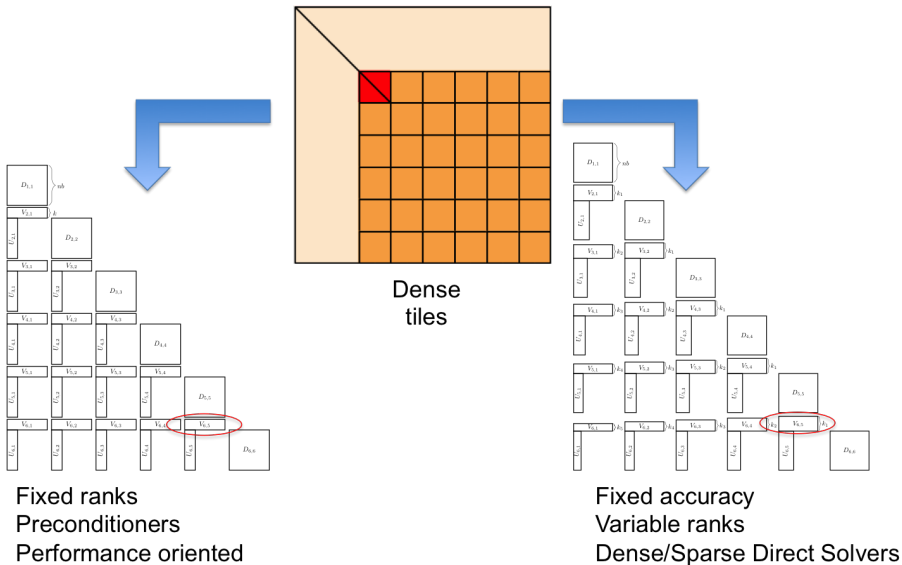
Outline

- 1 Motivations
- 2 Real Scientific Applications
- 3 Cholesky-based Matrix Computations
- 4 Climate/Weather Prediction Application
- 5 KBLAS

Computational Statistics for Climate/Weather Prediction Applications

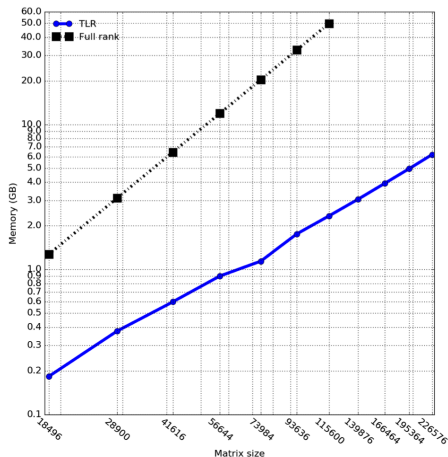
- Applications from climate and weather science often deal with a very large number of measurements regularly or irregularly located in geographical region.
- In geospatial statistics, these data are usually modeled as a realization from Gaussian spatial random field.
- This translates into evaluating the log-likelihood function, involving a large dense (but data-sparse) covariance matrix.

Dense Linear Algebra Renaissance



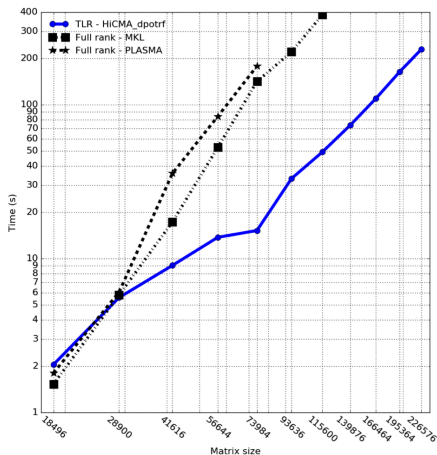
Tile Low Rank Cholesky: Memory Footprint

acc =
1e-9



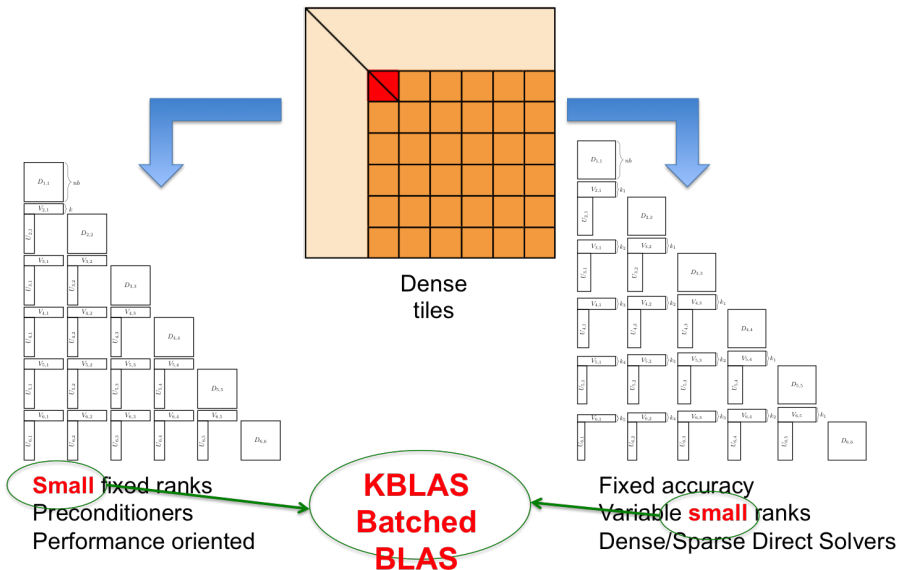
Akbudak et al., accepted at ISC'17

Tile Low Rank Cholesky: Time to Solution

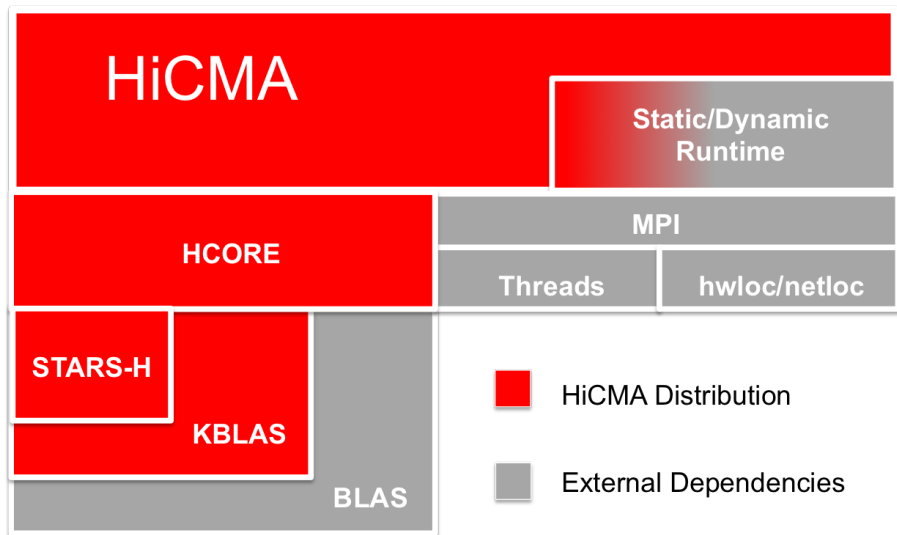


Akbudak et al., accepted at ISC'17

Dense Linear Algebra Renaissance



HiCMA Software Stack



Outline

- 1 Motivations
- 2 Real Scientific Applications
- 3 Cholesky-based Matrix Computations
- 4 Climate/Weather Prediction Application
- 5 KBLAS

Recursive formulation

- Usually used for Level 2 BLAS algorithms (e.g., panel factorization)
- Increase data locality
- Run at the cache level speed
- Again, not new and literature is quite rich: Kågström et. al (1998), Goto et. al (2008), etc.
- And it does pay off for Level 3 BLAS too!

Triangular matrix-matrix multiplication (TRMM)

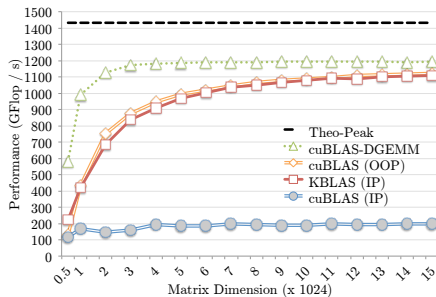


Figure: Performance comparisons of KBLAS DTRMM against IP and OOP cuBLAS DTRMM running on NVIDIA K40 GPU.

A. Charara, H. Ltaief and D. Keyes, Best Papers, EuroPar, 2016.
Integrated in CUDA 8.0

Triangular Solves (TRSM)

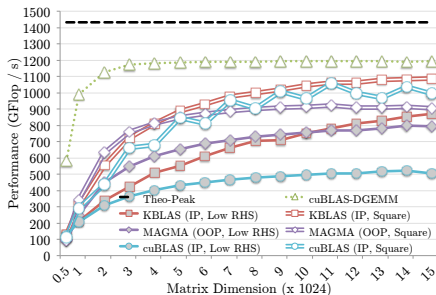


Figure: Performance comparisons of KBLAS IP DTRSM against cuBLAS IP DTRSM and MAGMA OOP TRSM running on NVIDIA K40 GPU, with square and low RHS matrices.

A. Charara, H. Ltaief and D. Keyes, Best Papers, EuroPar, 2016.
Integrated in CUDA 8.0

Advanced Batched BLAS Operations: HBLAS

Context:

- **Very** small sizes!
- Batch operation executions at each level of the tree
- Currently fixed sizes (need to handle variable sizes)
- Recursive formulation, stressing register usage
- Convert into batch of large GEMMs
- Minimize data transfer
- Enhance data locality
- Increase arithmetic intensity
- State-of-the-art implementations not well optimized for this scope or not supported

Advanced Batched BLAS Operations: HBLAS

HBLAS Matrix computations:

- Level 3 BLAS: SYRK, TRMM, TRSM
- Factorizations: POTRF
- Solves: POTRS, POSV, POTRI, POTI

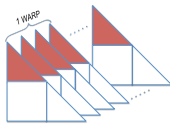
HBLAS Matrix compression:

- Batch QR factorizations
- Batch SVD

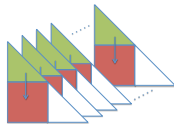
Advanced Batched BLAS Operations: HBLAS

Batches of Batched

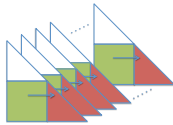
- Rec. Batch_DPOTRF



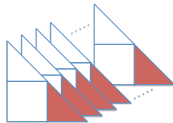
- Rec. Batch_DTRSM



- Rec. Batch_DSYRK



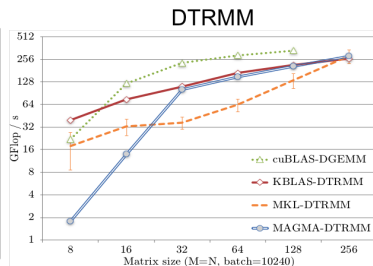
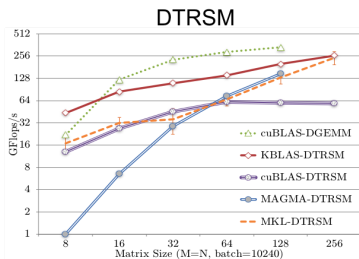
- Rec. Batch_DPOTRF



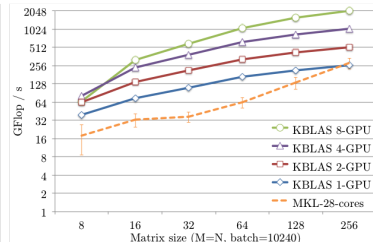
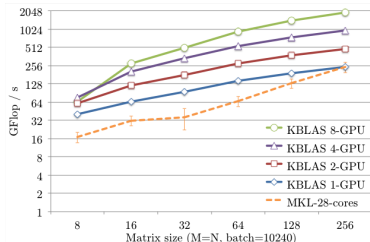
Profiling shows 76% of time is spent in batch DGEMM (MAGMABLAS).

Performance Results: Batched Level 3 BLAS on NVIDIA K40 GPUs

Single GPU



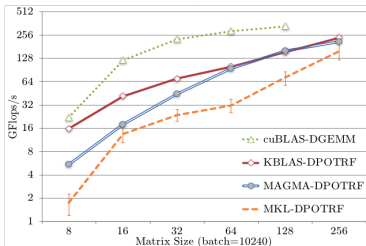
Multiple GPUs



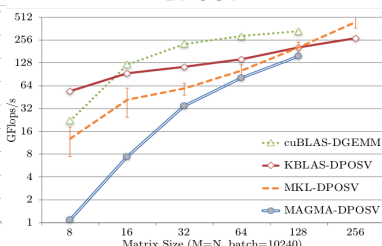
Performance Results: Batched Solves on NVIDIA K40 GPUs

Single GPU

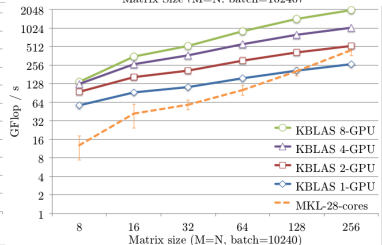
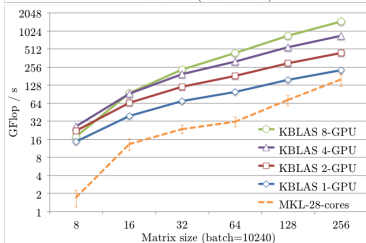
DPOTRF



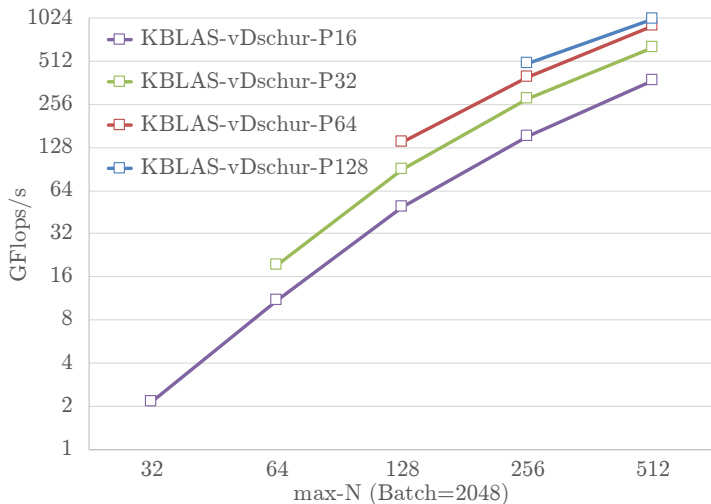
DPOSV



Multiple GPUs



Performance Results: Batched Schur Complement on NVIDIA K40 GPUs



Students/Collaborators/Support

- Extreme Computing Research Center @ KAUST: **W. Boukaram, A. Charara, G. Chávez, D. Keyes, D. Sukkari and G. Turkiyyah**
- L'Observatoire de Paris, LESIA: **R. Dembet, N. Doucet, E. Gendron, D. Gratadour, C. Morely, A. Sevin and F. Vidal**
- Innovative Computing Laboratory @ UTK: **PLASMA/MAGMA/ParSEC Teams**
- Barcelona Supercomputing Center, Spain: **R. Badia, P. Bellens, J. Labarta, X. Martorell, G. Miranda and S. Zhuang**
- Tokyo Institute of Technology: **R. Yokota**
- INRIA/INP Bordeaux, France: **E. Agullo, M. Faverge, F. Pruvost, M. Sergent and S. Thibault**
- KAUST Supercomputing Lab and IT Research Computing support
- NVIDIA GPU Research Center
- Intel Parallel Computing Center
- Cray Center of Excellence

