### Cholesky Factorization on Tile Low-Rank Matrices for Distributed-Memory Systems

## Kadir Akbudak, Hatem Ltaief, Aleksandr Mikhalev, Ali Charara and David keyes

#### Extreme Computing Research Center King Abdullah University of Science and Technology

IXPUG Workshop, KAUST

April 23-25, 2018



### How does future of Dense Linear Algebra (DLA) look like?

• Apparently dense matrices arising in scientific applications.

- Apparently dense matrices arising in scientific applications.
- Dense matrices might be compressed.

- Apparently dense matrices arising in scientific applications.
- Dense matrices might be compressed.
  - Tile low rank (TLR) matrix format



- Apparently dense matrices arising in scientific applications.
- Dense matrices might be compressed.
  - Tile low rank (TLR) matrix format
  - Cholesky factorization (for distributed-memory architectures)



- Apparently dense matrices arising in scientific applications.
- Dense matrices might be compressed.
  - Tile low rank (TLR) matrix format
  - Cholesky factorization (for distributed-memory architectures)
  - Huge performance improvement via cutting down flops



- Apparently dense matrices arising in scientific applications.
- Dense matrices might be compressed.
  - Tile low rank (TLR) matrix format
  - Cholesky factorization (for distributed-memory architectures)
  - Huge performance improvement via cutting down flops
  - Significantly less memory via storing small  $U_{i,j}$ and  $V_{i,j}$  matrices instead of big  $A_{i,j}$ , where  $A_{i,} = U_{i,j}V_{i,j}$



- Apparently dense matrices arising in scientific applications.
- Dense matrices might be compressed.
  - Tile low rank (TLR) matrix format
  - Cholesky factorization (for distributed-memory architectures)
  - Huge performance improvement via cutting down flops
  - Significantly less memory via storing small U<sub>i,j</sub> and V<sub>i,j</sub> matrices instead of big A<sub>i,j</sub>, where A<sub>i</sub> = U<sub>i,j</sub>V<sub>i,j</sub>
  - Preserving the accuracy requirements of the scientific application



- Apparently dense matrices arising in scientific applications.
- Dense matrices might be compressed.
  - Tile low rank (TLR) matrix format
  - Cholesky factorization (for distributed-memory architectures)
  - Huge performance improvement via cutting down flops
  - Significantly less memory via storing small U<sub>i,j</sub> and V<sub>i,j</sub> matrices instead of big A<sub>i,j</sub>, where A<sub>i</sub> = U<sub>i,j</sub>V<sub>i,j</sub>
  - Preserving the accuracy requirements of the scientific application
- MKL ScaLAPACK



- Apparently dense matrices arising in scientific applications.
- Dense matrices might be compressed.
  - Tile low rank (TLR) matrix format
  - Cholesky factorization (for distributed-memory architectures)
  - Huge performance improvement via cutting down flops
  - Significantly less memory via storing small U<sub>i,j</sub> and V<sub>i,j</sub> matrices instead of big A<sub>i,j</sub>, where A<sub>i</sub> = U<sub>i,j</sub>V<sub>i,j</sub>
  - Preserving the accuracy requirements of the scientific application
- MKL ScaLAPACK
  - pdpotrf → pdpotrf\_tlr?



### Climate/Weather Forecasting

• Computational statistics: multivariate large spatial data sets in climate/weather modeling:

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

(a) Problem Definition.



K. Akbudak et al.

w/Y. Sun and M. Genton  $\frac{3}{55}$ 

### Earth Science

• Seismic imaging: Imaging the subsurface by solving the Helmholtz equation (acoustic wave equation):

$$(-\bigtriangleup - k^2)u(x,w) = s(x,w)$$

 $k = \frac{w}{v(x)}$ , w is the angular frequency, v(x) is the seismic velocity field, and u(x, w) is the time-harmonic wavefield solution to the forcing term s(x, w).





w/ S. Zampini

### Materials Science

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



(b) Electronic structure.

Figure: Design of new materials.

w/ U. Schwingenschlogl

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

(a) Problem Definition.

• Symmetric, positive-definite matrix

- Symmetric, positive-definite matrix
- (Apparently) dense matrices

- Symmetric, positive-definite matrix
- (Apparently) dense matrices
- Often data-sparse

- Symmetric, positive-definite matrix
- (Apparently) dense matrices
- Often data-sparse
  - Decay of parameter correlations with distance

### Common Kernel: Cholesky Factorization (potrf)

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.

#### Ranks after Compression into Tile Low Rank (TLR) Format

• TLR format with varying ranks



#### Ranks after Compression into Tile Low Rank (TLR) Format

• TLR format with varying ranks



 Geospatial statistics, matrix size: 20000, block size: 500, accuracy threshold: 10<sup>-9</sup>, 2D problem

0	-500	133	47	35	154	83	44	33	59	49	38	30	40	37	33	29	33	32	30	27-
	133	500	139	48	86	163	86	44	50	59	50	38	37	41	37	33	32	33	32	30
	47	139	500	137	44	86	153	83	38	49	59	49	33	37	41	37	30	32	33	32
	35	48	137	500	33	44	83	164	30	38	49	59	29	33	37	41	27	30	32	33
	154	86	44	33	500	134	48	34	165	84	44	33	59	50	38	30	41	37	33	30
5	- 83	163	86	44	134	500	139	48	86	163	85	44	49	59	50	38	37	40	37	33-
	44	86	153	83	48	139	500	137	44	86	172	86	38	50	59	49	33	37	40	37
	33	44	83	164	34	48	137	500	33	44	86	166	30	39	50	59	29	33	37	41
	59	50	38	30	165	86	44	33	500	143	48	35	164	85	44	33	59	49	38	31
	49	59	49	38	84	163	86	44	143	500	143	48	84	159	87	44	49	59	49	38
10	- 38	50	59	49	44	85	172	86	48	143	500	134	44	86	156	81	38	49	58	49 -
	30	38	49	59	33	44	86	166	35	48	134	500	33	45	86	157	30	39	49	59
	40	37	33	29	59	49	38	30	164	84	44	33	500	138	48	35	162	86	44	33
	37	41	37	33	50	59	50	39	85	159	86	45	138	500	142	48	85	165	85	45
	33	37	41	37	38	50	59	50	44	87	156	86	48	142	500	133	44	84	159	85
15	- 29	33	37	41	30	38	49	59	33	44	81	157	35	48	133	500	33	44	81	157-
	33	32	30	27	41	37	33	29	59	49	38	30	162	85	44	33	500	142	47	34
	32	33	32	30	37	40	37	33	49	59	49	39	86	165	84	44	142	500	136	48
	30	32	33	32	33	37	40	37	38	49	58	49	44	85	159	81	47	136	500	130
	27	30	32	33	30	33	37	41	31	38	49	59	33	45	85	157	34	48	130	500
						5					10	-				15				

### Maximum Ranks after Compression into Tile Low Rank (TLR) Format

• Seismic imaging: 3D Helmholtz with N = 128 and  $k = N \times 0.625$  (10 grid points / wavelength) on  $\Omega = [0, 1]^3$ 



۲

- Seismic imaging: 3D Helmholtz with N = 128 and  $k = N \times 0.625$  (10 grid points / wavelength) on  $\Omega = [0, 1]^3$
- Heat map
- x axis: different tile sizes



- Seismic imaging: 3D Helmholtz with N = 128 and  $k = N \times 0.625$  (10 grid points / wavelength) on  $\Omega = [0, 1]^3$
- Heat map
- x axis: different tile sizes
- y axis: different accuracy thresholds



- Seismic imaging: 3D Helmholtz with N = 128 and  $k = N \times 0.625$  (10 grid points / wavelength) on  $\Omega = [0, 1]^3$
- Heat map 800 1e-14 828 • x axis: different 700 tile sizes 1e-12 • y axis: different <sup>600</sup> Maximum observed <sup>500</sup> 400 bserved threshold 1e-10 accuracy thresholds 1e-08 value/color: • value/color: 1e-06 TLR matrix for a 200 rank 0.0001 specific tile size and accuracy 100 0.01 256 512 1024 2048

#### Maximum Ranks after Compression into Tile Low Rank (TLR) Format

• Seismic imaging: 3D Helmholtz with N = 128 and  $k = N \times 0.625$  (10 grid points / wavelength) on  $\Omega = [0, 1]^3$ 



Tile size

K. Akbudak et al.

### Maximum Ranks after Compression into Tile Low Rank (TLR) Format

- Seismic imaging: 3D Helmholtz with N = 128 and  $k = N \times 0.625$  (10 grid points / wavelength) on  $\Omega = [0, 1]^3$
- Heat map
  Even for accuracy of 1e-14, ranks are smaller than

half of tile size.



- Seismic imaging: 3D Helmholtz with N = 128 and  $k = N \times 0.625$  (10 grid points / wavelength) on  $\Omega = [0, 1]^3$
- Heat map
- Even for accuracy of 1e-14, ranks are smaller than half of tile size.
- Even for tiles of size 2048, rank is smaller than 2048/2=1024.



#### Maximum Ranks after Compression into Tile Low Rank (TLR) Format

• Electronic structure: Hamiltonian/Overlap matrix

Heat map



Tile size

#### Maximum Ranks after Compression into Tile Low Rank (TLR) Format

• Electronic structure: Hamiltonian/Overlap matrix



- Electronic structure: Hamiltonian/Overlap matrix
- Heat map
- x axis: different tile sizes
- y axis: different accuracy thresholds



#### Maximum Ranks after Compression into Tile Low Rank (TLR) Format

• Electronic structure: Hamiltonian/Overlap matrix



Tile size

#### Maximum Ranks after Compression into Tile Low Rank (TLR) Format

• Electronic structure: Hamiltonian/Overlap matrix



Tile size

- Electronic structure: Hamiltonian/Overlap matrix
- Heat map
  Even for accuracy of 1e-14, maximum rank is 55 which is considerably smaller than tile size.



Tile size

- Electronic structure: Hamiltonian/Overlap matrix
- Heat map
  Even for accuracy of 1e-14, maximum rank is 55 which is considerably smaller than tile size.
- Even for larger tiles, maximum rank is 55.



Tile size

# The HiCMA Library: Hierarchical Computations on Manycore Architectures



Available at http://github.com/ecrc/hicma
# Dense Linear Algebra Renaissance



The tile low-rank (TLR) Cholesky algorithm can be expressed with the following four computational kernels:

hcore\_dpotrf: The kernel performs the Cholesky factorization of a diagonal (lower triangular) tile. It is similar to DPOTRF since the diagonal tiles are dense.

The tile low-rank (TLR) Cholesky algorithm can be expressed with the following four computational kernels:

- hcore\_dpotrf: The kernel performs the Cholesky factorization of a diagonal (lower triangular) tile. It is similar to DPOTRF since the diagonal tiles are dense.
- hcore\_dtrsm: The operation applies an update to an off-diagonal low-rank tile of the input matrix, resulting from factorization of the diagonal tile above it and overrides it with the final elements of the output matrix:  $V_{(i,k)} = V_{(i,k)} \times D_{(k,k)}^{-1}$ . The operation is a triangular solve.

The tile low-rank (TLR) Cholesky algorithm can be expressed with the following four computational kernels:

- hcore\_dpotrf: The kernel performs the Cholesky factorization of a diagonal (lower triangular) tile. It is similar to DPOTRF since the diagonal tiles are dense.
- hcore\_dtrsm: The operation applies an update to an off-diagonal low-rank tile of the input matrix, resulting from factorization of the diagonal tile above it and overrides it with the final elements of the output matrix:  $V_{(i,k)} = V_{(i,k)} \times D_{(k,k)}^{-1}$ . The operation is a triangular solve.
- hcore\_dsyrk: The kernel applies updates to a diagonal (lower triangular) tile of the input matrix, resulting from factorization of the low-rank tiles to the left of it:  $D_{(j,j)} = D_{(j,j)} - (U_{(j,k)} \times V_{(j,k)}^T) \times (U_{(j,k)} \times V_{(j,k)}^T)^T$ . The operation is a symmetric rank-k update.

The tile low-rank (TLR) Cholesky algorithm can be expressed with the following four computational kernels:

- hcore\_dpotrf: The kernel performs the Cholesky factorization of a diagonal (lower triangular) tile. It is similar to DPOTRF since the diagonal tiles are dense.
  - hcore\_dtrsm: The operation applies an update to an off-diagonal low-rank tile of the input matrix, resulting from factorization of the diagonal tile above it and overrides it with the final elements of the output matrix:  $V_{(i,k)} = V_{(i,k)} \times D_{(k,k)}^{-1}$ . The operation is a triangular solve.
  - hcore\_dsyrk: The kernel applies updates to a diagonal (lower triangular) tile of the input matrix, resulting from factorization of the low-rank tiles to the left of it:  $D_{(j,j)} = D_{(j,j)} - (U_{(j,k)} \times V_{(j,k)}^T) \times (U_{(j,k)} \times V_{(j,k)}^T)^T$ . The operation is a symmetric rank-*k* update.
  - hcore\_dgemm: The operation applies updates to an off-diagonal low-rank tile of the input matrix, resulting from factorization of the low-rank tiles to the left of it. The operation involves two QR factorizations, one reduced SVD (depending on the rank and/or the accuracy parameter) and two matrix-matrix multiplications.

#### **Algorithm 1** hicma\_dpotrf (HicmaLower, D, U, V, N, nb, rank, acc) p = N / nb $D_{1,1}$ for k = 1 to p do $V_{2.1}$ #task access(read&write:D(k)) $hcore_dpotrf$ (HicmaLower, D(k)) $D_{2,2}$ for i = k+1 to p do #task access(read:D(k), read&write:V(i,k)) $D_{3.3}$ $hcore_dtrsm(V(i,k), D(k,k))$ end for for i = k+1 to p do # task access(read:U(j,k), read:V(j,k), read&write:D(j)) hcore\_dsyrk (D(j), U(j,k), V(j,k)) $D_{5.5}$ for i = j+1 to p do #task access(read:U(i,k), Ves read:V(i,k)), read:U(j,k), read:V(j,k), $D_{6.6}$ read&write:U(i,j), read&write:V(i,j)) hcore\_dgemm (U(i,k), V(i,k), U(j,k), V(j,k), U(i,j), V(i,j), rank, acc) end for end for end for



- Execution times for MKL dpotrf (dense) vs HiCMA dpotrf (TLR). Smaller the better.
- Application: Geospatial statistic, square exp. kernel



- Execution times for MKL dpotrf (dense) vs HiCMA dpotrf (TLR). Smaller the better.
- Application: Geospatial statistic, square exp. kernel
- Accuracy threshold:  $10^{-9}$ . This accuracy level is enough for a large amount of cases encountered in the geospatial statistics application.



- Execution times for MKL dpotrf (dense) vs HiCMA dpotrf (TLR). Smaller the better.
- Application: Geospatial statistic, square exp. kernel
- Accuracy threshold:  $10^{-9}$ . This accuracy level is enough for a large amount of cases encountered in the geospatial statistics application.
- StarPU (task-based runtime)



- Execution times for MKL dpotrf (dense) vs HiCMA dpotrf (TLR). Smaller the better.
- Application: Geospatial statistic, square exp. kernel
- Accuracy threshold: 10<sup>-9</sup>. This accuracy level is enough for a large amount of cases encountered in the geospatial statistics application.
- StarPU (task-based runtime)
- Three architectures:
  - Sandy Bridge (SNB): AVX, 2.0GHz
     2 sockets x 8 cores
  - Haswell (HSW): AVX2, 2.4GHz
    2 sockets x 18 cores
  - Skylake (SKL): AVX512, 2.1GHz 2 sockets × 28 cores



- Execution times for MKL dpotrf (dense) vs HiCMA dpotrf (TLR). Smaller the better.
- x axis (log-scale): different matrix sizes



- Execution times for MKL dpotrf (dense) vs HiCMA dpotrf (TLR). Smaller the better.
- x axis (*log*-scale): different matrix sizes
  - Missing points due to not enough memory for larger matrices



- Execution times for MKL dpotrf (dense) vs HiCMA dpotrf (TLR). Smaller the better.
- x axis (*log*-scale): different matrix sizes
  - Missing points due to not enough memory for larger matrices

K. Akbudak et al.

• y axis (*log*-scale): time of dpotrf in seconds



- Execution times for MKL dpotrf (dense) vs HiCMA dpotrf (TLR). Smaller the better.
- x axis (*log*-scale): different matrix sizes
  - Missing points due to not enough memory for larger matrices
- y axis (*log*-scale): time of dpotrf in seconds
- Smaller time as chips get better



- Execution times for MKL dpotrf (dense) vs HiCMA dpotrf (TLR). Smaller the better.
- x axis (*log*-scale): different matrix sizes
  - Missing points due to not enough memory for larger matrices
- y axis (*log*-scale): time of dpotrf in seconds
- Smaller time as chips get better
- Solve larger problems



- Execution times for MKL dpotrf (dense) vs HiCMA dpotrf (TLR). Smaller the better.
- x axis (*log*-scale): different matrix sizes
  - Missing points due to not enough memory for larger matrices
- y axis (*log*-scale): time of dpotrf in seconds
- Smaller time as chips get better
- Solve larger problems
- Smaller slope:  $O(n^3) \rightarrow O(kn^2)$



Experimental Results

#### ScaLAPACK vs TLR Cholesky (Shaheen-2, HSW, acc=10<sup>-9</sup>)

• Execution times for ScaLAPACK dpotrf (dense) vs HiCMA dpotrf (TLR). Smaller the better.



K. Akbudak, H. Ltaief, A. Mikhalev, A. Charara, and D. E. Keyes, Exploiting Data Sparsity for Large-Scale Matrix Computations, Submitted to EuroPar, 2018.

- Execution times for ScaLAPACK dpotrf (dense) vs HiCMA dpotrf (TLR). Smaller the better.
- Application: Geospatial statistic, square exp. kernel



K. Akbudak, H. Ltaief, A. Mikhalev, A. Charara, and D. E. Keyes, Exploiting Data Sparsity for Large-Scale Matrix Computations, Submitted to EuroPar, 2018.

- Execution times for ScaLAPACK dpotrf (dense) vs HiCMA dpotrf (TLR). Smaller the better.
- Application: Geospatial statistic, square exp. kernel
- Accuracy threshold:  $10^{-9}$ .



K. Akbudak, H. Ltaief, A. Mikhalev, A. Charara, and D. E. Keyes, Exploiting Data Sparsity for Large-Scale Matrix Computations, Submitted to EuroPar, 2018.

- Execution times for ScaLAPACK dpotrf (dense) vs HiCMA dpotrf (TLR). Smaller the better.
- Application: Geospatial statistic, square exp. kernel
- Accuracy threshold:  $10^{-9}$ .
- StarPU (task-based runtime)



K. Akbudak, H. Ltaief, A. Mikhalev, A. Charara, and D. E. Keyes, Exploiting Data Sparsity for Large-Scale Matrix Computations, Submitted to EuroPar, 2018.

- Execution times for ScaLAPACK dpotrf (dense) vs HiCMA dpotrf (TLR). Smaller the better.
- Application: Geospatial statistic, square exp. kernel
- Accuracy threshold:  $10^{-9}$ .
- StarPU (task-based runtime)
- x axis (*log*-scale): different matrix sizes



K. Akbudak, H. Ltaief, A. Mikhalev, A. Charara, and D. E. Keyes, Exploiting Data Sparsity for Large-Scale Matrix Computations, Submitted to EuroPar, 2018.

- Execution times for ScaLAPACK dpotrf (dense) vs HiCMA dpotrf (TLR). Smaller the better.
- Application: Geospatial statistic, square exp. kernel
- Accuracy threshold:  $10^{-9}$ .
- StarPU (task-based runtime)
- x axis (*log*-scale): different matrix sizes
- y axis (*log*-scale): time of dpotrf in seconds



K. Akbudak, H. Ltaief, A. Mikhalev, A. Charara, and D. E. Keyes, Exploiting Data Sparsity for Large-Scale Matrix Computations, Submitted to EuroPar, 2018.

- Execution times for ScaLAPACK dpotrf (dense) vs HiCMA dpotrf (TLR). Smaller the better.
- Application: Geospatial statistic, square exp. kernel
- Accuracy threshold:  $10^{-9}$ .
- StarPU (task-based runtime)
- x axis (*log*-scale): different matrix sizes
- y axis (*log*-scale): time of dpotrf in seconds
- 85x speedup



K. Akbudak, H. Ltaief, A. Mikhalev, A. Charara, and D. E. Keyes, Exploiting Data Sparsity for Large-Scale Matrix Computations, Submitted to EuroPar, 2018.

- Execution times for ScaLAPACK dpotrf (dense) vs HiCMA dpotrf (TLR). Smaller the better.
- Application: Geospatial statistic, square exp. kernel
- Accuracy threshold:  $10^{-9}$ .
- StarPU (task-based runtime)
- x axis (*log*-scale): different matrix sizes
- y axis (*log*-scale): time of dpotrf in seconds



K. Akbudak, H. Ltaief, A. Mikhalev, A. Charara, and D. E. Keyes, Exploiting Data Sparsity for Large-Scale Matrix Computations, Submitted to EuroPar, 2018. Experimental Results

# TLR Cholesky: Memory Footprint (acc= $10^{-9}$ )

 Memory footprint for dense vs TLR matrices.
 Smaller the better.



# TLR Cholesky: Memory Footprint (acc= $10^{-9}$ )

- Memory footprint for dense vs TLR matrices. Smaller the better.
- Application: Geospatial statistic, square exp. kernel



# TLR Cholesky: Memory Footprint (acc= $10^{-9}$ )

- Memory footprint for dense vs TLR matrices. Smaller the better.
- Application: Geospatial statistic, square exp. kernel
- Accuracy threshold:  $10^{-9}$



- Memory footprint for dense vs TLR matrices. Smaller the better.
- Application: Geospatial statistic, square exp. kernel
- Accuracy threshold:  $10^{-9}$
- x axis (*log*-scale): different matrix sizes (*N*)



- Memory footprint for dense vs TLR matrices. Smaller the better.
- Application: Geospatial statistic, square exp. kernel
- Accuracy threshold:  $10^{-9}$
- x axis (*log*-scale): different matrix sizes (*N*)
- y axis (*log*-scale): memory required for storing matrix in terms of giga bytes



- Memory footprint for dense vs TLR matrices. Smaller the better.
- Application: Geospatial statistic, square exp. kernel
- Accuracy threshold:  $10^{-9}$
- x axis (*log*-scale): different matrix sizes (*N*)
- y axis (*log*-scale): memory required for storing matrix in terms of giga bytes
- Full rank: N<sup>2</sup>



- Memory footprint for dense vs TLR matrices. Smaller the better.
- Application: Geospatial statistic, square exp. kernel
- Accuracy threshold: 10<sup>-9</sup>
- x axis (*log*-scale): different matrix sizes (*N*)
- y axis (*log*-scale): memory required for storing matrix in terms of giga bytes
- Full rank: N<sup>2</sup>
- TLR:  $\frac{N^2k}{nb}$



- Memory footprint for dense vs TLR matrices. Smaller the better.
- Application: Geospatial statistic, square exp. kernel
- Accuracy threshold:  $10^{-9}$
- x axis (*log*-scale): different matrix sizes (*N*)
- y axis (*log*-scale): memory required for storing matrix in terms of giga bytes
- Full rank: N<sup>2</sup>
- TLR:  $\frac{N^2k}{nb}$

• k: average rank



- Memory footprint for dense vs TLR matrices. Smaller the better.
- Application: Geospatial statistic, square exp. kernel
- Accuracy threshold: 10<sup>-9</sup>
- x axis (*log*-scale): different matrix sizes (*N*)
- y axis (*log*-scale): memory required for storing matrix in terms of giga bytes
- Full rank: N<sup>2</sup>
- TLR:  $\frac{N^2k}{nb}$ 
  - k: average rank
  - nb: tile size



#### Impact of Accuracy Thresholds, 1M Matrix Size, nb=2700



K. Akbudak, H. Ltaief, A. Mikhalev, A. Charara, and D. E. Keyes, Exploiting Data Sparsity for Large-Scale Matrix Computations, Submitted to EuroPar, 2018.

#### Impact of Accuracy Thresholds, 1M Matrix Size, nb=2700



K. Akbudak, H. Ltaief, A. Mikhalev, A. Charara, and D. E. Keyes, Exploiting Data Sparsity for Large-Scale Matrix Computations, Submitted to EuroPar, 2018.
#### Impact of Accuracy Thresholds, 1M Matrix Size, nb=2700

- Memory footprint for dense vs TLR matrices. Smaller the better.
- Application: Geospatial statistic, square exp. kernel
- x axis (*log*-scale): different accuracy thresholds



K. Akbudak, H. Ltaief, A. Mikhalev, A. Charara, and D. E. Keyes, Exploiting Data Sparsity for Large-Scale Matrix Computations, Submitted to EuroPar, 2018.

K. Akbudak et al.

#### Impact of Accuracy Thresholds, 1M Matrix Size, nb=2700

- Memory footprint for dense vs TLR matrices. Smaller the better.
- Application: Geospatial statistic, square exp. kernel
- x axis (*log*-scale): different accuracy thresholds
- y axis (*log*-scale): memory required for storing matrix in terms of giga bytes



K. Akbudak, H. Ltaief, A. Mikhalev, A. Charara, and D. E. Keyes, Exploiting Data Sparsity for Large-Scale Matrix Computations, Submitted to EuroPar, 2018.

#### Impact of Accuracy Thresholds, 1M Matrix Size, nb=2700

- Memory footprint for dense vs TLR matrices. Smaller the better.
- Application: Geospatial statistic, square exp. kernel
- x axis (*log*-scale): different accuracy thresholds
- y axis (*log*-scale): memory required for storing matrix in terms of giga bytes
- Being more accurate requires higher ranks so larger memory is used.



K. Akbudak, H. Ltaief, A. Mikhalev, A. Charara, and D. E. Keyes, Exploiting Data Sparsity for Large-Scale Matrix Computations, Submitted to EuroPar, 2018.

K. Akbudak et al.

• Execution times for HiCMA dpotrf (TLR)



- Execution times for HiCMA dpotrf (TLR)
- Application: Geospatial statistic, square exp. kernel



- Execution times for HiCMA dpotrf (TLR)
- Application: Geospatial statistic, square exp. kernel
- x axis (*log*-scale): different accuracy thresholds



K. Akbudak, H. Ltaief, A. Mikhalev, A. Charara, and D. E. Keyes, Exploiting Data Sparsity for Large-Scale Matrix Computations, Submitted to EuroPar, 2018.

- Execution times for HiCMA dpotrf (TLR)
- Application: Geospatial statistic, square exp. kernel
- x axis (*log*-scale): different accuracy thresholds
- y axis (*log*-scale): time of dpotrf in seconds



K. Akbudak, H. Ltaief, A. Mikhalev, A. Charara, and D. E. Keyes, Exploiting Data Sparsity for Large-Scale Matrix Computations, Submitted to EuroPar, 2018.

- Execution times for HiCMA dpotrf (TLR)
- Application: Geospatial statistic, square exp. kernel
- x axis (*log*-scale): different accuracy thresholds
- y axis (*log*-scale): time of dpotrf in seconds
- Being more accurate requires higher ranks so more flops are done.





K. Akbudak et al.



Computations, Submitted to EuroPar, 2018.

- Execution times for HiCMA dpotrf (TLR) for different number of nodes
- Application: Geospatial statistic, square exp. kernel
- x axis (*log*-scale): different matrix sizes



- Execution times for HiCMA dpotrf (TLR) for different number of nodes
- Application: Geospatial statistic, square exp. kernel
- x axis (*log*-scale): different matrix sizes
- y axis (*log*-scale): time of dpotrf in minutes



K. Akbudak, H. Ltaief, A. Mikhalev, A. Charara, and D. E. Keyes, Exploiting Data Sparsity for Large-Scale Matrix Computations, Submitted to EuroPar, 2018.

- Execution times for HiCMA dpotrf (TLR) for different number of nodes
- Application: Geospatial statistic, square exp. kernel
- x axis (*log*-scale): different matrix sizes
- y axis (*log*-scale): time of dpotrf in minutes
- Running ScaLAPACK for larger cases is not feasible because of time and memory.



- Execution times for HiCMA dpotrf (TLR) for different number of nodes
- Application: Geospatial statistic, square exp. kernel
- x axis (*log*-scale): different matrix sizes
- y axis (*log*-scale): time of dpotrf in minutes
- Running ScaLAPACK for larger cases is not feasible because of time and memory.
- Missing points for larger matrices: not enough memory



- Execution times for HiCMA dpotrf (TLR) for different number of nodes
- Application: Geospatial statistic, square exp. kernel
- x axis (*log*-scale): different matrix sizes
- y axis (*log*-scale): time of dpotrf in minutes
- Running ScaLAPACK for larger cases is not feasible because of time and memory.
- Missing points for larger matrices: not enough memory
- Missing points for smaller matrices: not enough work



• Execution times for HiCMA dpotrf (TLR) for different number of nodes



K. Akbudak, H. Ltaief, A. Mikhalev, A. Charara, and D. E. Keyes, Exploiting Data Sparsity for Large-Scale Matrix Computations, Submitted to EuroPar, 2018.

- Execution times for HiCMA dpotrf (TLR) for different number of nodes
- Application: Geospatial statistic, square exp. kernel



K. Akbudak, H. Ltaief, A. Mikhalev, A. Charara, and D. E. Keyes, Exploiting Data Sparsity for Large-Scale Matrix Computations, Submitted to EuroPar, 2018.

- Execution times for HiCMA dpotrf (TLR) for different number of nodes
- Application: Geospatial statistic, square exp. kernel
- x axis (*log*-scale): different matrix sizes



K. Akbudak, H. Ltaief, A. Mikhalev, A. Charara, and D. E. Keyes, Exploiting Data Sparsity for Large-Scale Matrix Computations, Submitted to EuroPar, 2018.

- Execution times for HiCMA dpotrf (TLR) for different number of nodes
- Application: Geospatial statistic, square exp. kernel
- x axis (*log*-scale): different matrix sizes
- y axis (log-scale): time of dpotrf in minutes



K. Akbudak, H. Ltaief, A. Mikhalev, A. Charara, and D. E. Keyes, Exploiting Data Sparsity for Large-Scale Matrix Computations, Submitted to EuroPar, 2018.

- Execution times for HiCMA dpotrf (TLR) for different number of nodes
- Application: Geospatial statistic, square exp. kernel
- x axis (*log*-scale): different matrix sizes
- y axis (log-scale): time of dpotrf in minutes
- We ran HiCMA on the latest Intel deployment.



K. Akbudak, H. Ltaief, A. Mikhalev, A. Charara, and D. E. Keyes, Exploiting Data Sparsity for Large-Scale Matrix Computations, Submitted to EuroPar, 2018.

Execution times for HiCMA dpotrf (TLR) for different number of nodes



Computations, Submitted to EuroPar, 2018.

- Execution times for HiCMA dpotrf (TLR) for different number of nodes
- Application: Geospatial statistic, square exp. kernel



K. Akbudak, H. Ltaief, A. Mikhalev, A. Charara, and D. E. Keyes, Exploiting Data Sparsity for Large-Scale Matrix Computations, Submitted to EuroPar, 2018.

- Execution times for HiCMA dpotrf (TLR) for different number of nodes
- Application: Geospatial statistic, square exp. kernel
- x axis (*log*-scale): different matrix sizes



K. Akbudak, H. Ltaief, A. Mikhalev, A. Charara, and D. E. Keyes, Exploiting Data Sparsity for Large-Scale Matrix Computations, Submitted to EuroPar, 2018.

- Execution times for HiCMA dpotrf (TLR) for different number of nodes
- Application: Geospatial statistic, square exp. kernel
- x axis (*log*-scale): different matrix sizes
- y axis (log-scale): time of dpotrf in minutes



K. Akbudak, H. Ltaief, A. Mikhalev, A. Charara, and D. E. Keyes, Exploiting Data Sparsity for Large-Scale Matrix Computations, Submitted to EuroPar, 2018.

- Execution times for HiCMA dpotrf (TLR) for different number of nodes
- Application: Geospatial statistic, square exp. kernel
- x axis (*log*-scale): different matrix sizes
- y axis (log-scale): time of dpotrf in minutes
- Not big difference between turbo on and off due low arithmetic intensity



# Traces Chameleon: Dense dpotrf time=18.1s on 4 nodes of Shaheen-2 with a matrix size of 54K



K. Akbudak et al.

# Traces HiCMA: Data-sparse dpotrf time=1.8s on 4 nodes of Shaheen-2 with a matrix size of 54K



K. Akbudak et al.

## Conclusion & Future Work

- We cut down flops through compression.
- We use less memory.
- Accuracy requirement is satisfied.
- HiCMA is for both shared and distributed memory systems.
- HODLR/H (non-nested bases) data compression format.
- More matrix computation algorithms (LU/QR, eigenvalue solvers/SVD, matrix inversion, etc.)
- Support for more large-scale scientific applications.

## The Hourglass Revisited



## The Hourglass Revisited



## Students/Collaborators/Vendors

- Extreme Computing Research Center @ KAUST: R. Abdelkhalek, S. Abdullah, K. Akbudak, R. AlOmairy, A. Alonazy, T. Alturkestani, W. Boukaram, A. Charara, N. Doucet, M. Genton, E. A. Gonzalez Fisher, D. Keyes, H. Ltaief, A. Mikhalev, D. Sukkari and Y. Sun
- Tokyo Institute of Technology, Japan: R. Yokota
- Innovative Computing Laboratory @ UTK, USA: PLASMA/MAGMA/PaRSEC Teams
- INRIA/INP/LaBRI Bordeaux, France: Runtime/HiePACS Teams
- Max-Planck Institute@Leipzig, Germany: R. Kriemann
- American University of Beirut, Lebanon: G. Turkiyyah
- KAUST Supercomputing Lab
- Intel Parallel Computing Center
- Cray Center of Excellence





Acknowledgments

## Thank You!

## Questions?

K. Akbudak et al.