### Scalable Algorithms for Tensor Computations

Edgar Solomonik

L P. N A @CS@Illinois

Department of Computer Science University of Illinois at Urbana-Champaign

SIAM PP, Seattle WA

# Laboratory for Parallel Numerical Algorithms

#### Recent/ongoing research topics

- parallel matrix computations
  - QR factorization
  - triangular solve
  - eigenvalue problems
- tensor computations
  - tensor decomposition
  - sparse tensor kernels
  - tensor completion

#### simulation of quantum systems

- tensor networks
- quantum chemistry
- quantum circuits
- fast bilinear algorithms
  - convolution algorithms
  - tensor symmetry
  - fast matrix multiplication







http://lpna.cs.illinois.edu

## Outline

### Introduction

- 2 Matrix Factorizations
- 3 Tensor Computations
- 4 Software Libraries



## 3D algorithms for matrix computations

For Cholesky factorization with p processors, BSP (critical path) costs are

$$F = \Theta(n^3/p), \quad W = \Theta(n^2/\sqrt{cp}), \quad S = \Theta(\sqrt{cp})$$

using c matrix copies (processor grid is 2D for c = 1, 3D for  $c = p^{1/3}$ ).

Achieving similar costs for LU, QR, and the symmetric eigenvalue problem requires algorithmic changes.

square TRSM $\sqrt{1}$	rectangular TRSM $\sqrt{2}$
pairwise pivoting $\sqrt{3}$	tournament pivoting $\checkmark^4$
Givens on square $\checkmark^3$	Householder on rect. $\sqrt{5}$
eigenvalues only $\checkmark^5$	eigenvectors X
	square TRSM $\sqrt{1}$ pairwise pivoting $\sqrt{3}$ Givens on square $\sqrt{3}$ eigenvalues only $\sqrt{5}$

 $\sqrt{\text{means costs attained (synchronization within polylog factors)}}$ .

<sup>1</sup>B. Lipshitz, MS thesis 2013
<sup>2</sup>T. Wicky, E.S., T. Hoefler, IPDPS 2017
<sup>3</sup>A. Tiskin, FGCS 2007
<sup>4</sup>E.S., J. Demmel, EuroPar 2011
<sup>5</sup>E.S., G. Ballard, T. Hoefler, J. Demmel, SPAA 2017

LPNA

## New algorithms can circumvent lower bounds

For TRSM, we can achieve a lower synchronization/communication cost by performing triangular inversion on diagonal blocks





Tobias Wicky

- $\bullet$  decreases synchronization cost by  ${\cal O}(p^{2/3})$  on p processors with respect to known algorithms
- optimal communication for any number of right-hand sides
- MS thesis work by Tobias Wicky<sup>1</sup>

<sup>&</sup>lt;sup>1</sup>T. Wicky, E.S., T. Hoefler, IPDPS 2017

### Cholesky-QR2 for rectangular matrices

Cholesky-QR2<sup>1</sup> with 3D Cholesky gives a practical 3D QR algorithm<sup>2</sup>

- Compute  $oldsymbol{A} = \hat{oldsymbol{Q}} oldsymbol{R}$  using Cholesky  $oldsymbol{A}^T oldsymbol{A} = oldsymbol{R}^T oldsymbol{R}$
- Correct computed factorization by Cholesky-QR of  $\hat{Q}$
- $\bullet$  Attains full accuracy so long as  ${\rm cond}({\pmb A}) < 1/\sqrt{\epsilon_{\rm mach}}$





Edward Hutter

<sup>1</sup>T. Fukaya, Y. Nakatsukasa, Y. Yanagisawa, Y. Yamamoto, 2014

<sup>2</sup>E. Hutter, E.S., IPDPS 2019

LPNA

### Tridiagonalization

Reducing the symmetric matrix  $oldsymbol{A} \in \mathbb{R}^{n imes n}$  to a tridiagonal matrix

$$T = Q^T A Q$$

by an orthogonal similarity transformation is most costly in eigenvalue computation (SVD is similar).

• can be done by successive column QR factorizations

$$T = \underbrace{Q_n^T \cdots Q_1^T}_{Q^T} A \underbrace{Q_1 \cdots Q_n}_{Q}$$

- two-sided updates harder to manage than one-sided
- can use n/b QRs on panels of b columns to go to band-width b+1
- b = 1 gives direct tridiagonalization

### Multi-stage tridiagonalization

Writing the orthogonal transformation in Householder form, we get

$$\underbrace{(\boldsymbol{I} - \boldsymbol{U}\boldsymbol{T}\boldsymbol{U}^T)^T}_{\boldsymbol{Q}^T}\boldsymbol{A}\underbrace{(\boldsymbol{I} - \boldsymbol{U}\boldsymbol{T}\boldsymbol{U}^T)}_{\boldsymbol{Q}} = \boldsymbol{A} - \boldsymbol{U}\boldsymbol{V}^T - \boldsymbol{V}\boldsymbol{U}^T$$

where columns of  $\boldsymbol{U}$  contains Householder vectors and  $\boldsymbol{V}$  is

$$oldsymbol{V}^T = oldsymbol{T}oldsymbol{U}^T + rac{1}{2}oldsymbol{T}^Toldsymbol{U}^T \underbrace{oldsymbol{A}oldsymbol{U}}_{ ext{bottleneck}}oldsymbol{T}oldsymbol{U}^T$$

if b = 1, U is a column-vector, and AU is dominated by vertical communication cost (moving A between memory and cache)

• idea: reduce to banded matrix  $(b \gg 1)$  first<sup>1</sup>

<sup>&</sup>lt;sup>1</sup>T. Auckenthaler, H. Bungartz, T. Huckle, L. Krämer, B. Lang, P. Willems 2011

# Successive band reduction (SBR)

After reducing to a banded matrix, we need to transform the banded matrix to a tridiagonal one

- fewer nonzeros lead to lower computational cost,  $F = O(n^2 b/p)$
- however, transformations introduce fill/bulges
- bulges must be chased down the band<sup>1</sup>



 communication- and synchronization-efficient 1D SBR algorithm known for small band-width<sup>2</sup>

<sup>2</sup>G. Ballard, J. Demmel, N. Knight 2012

<sup>&</sup>lt;sup>1</sup>B. Lang 1993; C.H. Bischof, B. Lang, X. Sun 2000

### Communication-efficient eigenvalue computation

Previous work (start-of-the-art): two-stage tridiagonalization

• implemented in ELPA, can outperform ScaLAPACK<sup>1</sup>

• with  $n = n/\sqrt{p}$ , 1D SBR gives  $W = O(n^2/\sqrt{p})$ ,  $S = O(\sqrt{p}\log^2(p))$ 

New results<sup>3</sup>: many-stage tridiagonalization

- $\Theta(\log(p))$  intermediate band-widths to achieve  $W = O(n^2/p^{2/3})$
- communication-efficient rectangular QR with processor groups



#### • 3D SBR (each QR and matrix multiplication update parallelized)

- <sup>2</sup>G. Ballard, J. Demmel, N. Knight 2012
- <sup>3</sup>E.S., G. Ballard, J. Demmel, T. Hoefler 2017

LPNA

<sup>&</sup>lt;sup>1</sup>T. Auckenthaler, H. Bungartz, T. Huckle, L. Krämer, B. Lang, P. Willems 2011

### Symmetric eigensolver results summary

W	Q	S
$n^2/\sqrt{p}$	$n^3/p$	$n\log(p)$
$n^2/\sqrt{p}$	-	$n\log(p)$
$n^2/\sqrt{p}$	$n^2 \log(n) / \sqrt{p}$	$\sqrt{p}(\log^2(p) + \log(n))$
$n^2/p^{2/3}$	$n^2\log(p)/p^{2/3}$	$p^{2/3}\log^2 p$
	$\frac{W}{n^2/\sqrt{p}} \\ \frac{n^2/\sqrt{p}}{n^2/\sqrt{p}} \\ \frac{n^2/p^{2/3}}{n^2/p^{2/3}}$	$\begin{array}{c c c c c c c c c c c c c c c c c c c $

- costs are asymptotic (same computational cost F for eigenvalues)
- W horizontal (interprocessor) communication
- Q vertical (memory–cache) communication excluding  $W+F/\sqrt{H}$  where H is cache size
- S synchronization cost (number of supersteps)

## **Tensor** Decompositions

- Tensor of order N has N modes and dimensions  $s\times \dots \times s$
- Canonical polyadic (CP) tensor decomposition<sup>1</sup>



- Alternating least squares (ALS) is most widely used method
  - Monotonic linear convergence
- Gauss-Newton method is an emerging alternative
  - Non-monotonic, but can achieve superlinear convergence rate

<sup>&</sup>lt;sup>1</sup>T.G. Kolda and B.W. Bader, SIAM Review 2009

## Accelerating Alternating Least Squares



New algorithm: pairwise perturbation  $(PP)^1$  approximates ALS

- accurate when ALS sweep over factor matrices stagnates
- rank  $R < s^{N-1}$  CP decomposition:
  - ALS sweep cost  ${\cal O}(s^NR) \Rightarrow {\cal O}(s^2R),$  up to 600x speed-up
- rank R < s Tucker decomposition:
  - ALS sweep cost  $O(s^N R) \Rightarrow O(s^2 R^{N-1})$



Linjian Ma

<sup>&</sup>lt;sup>1</sup>L. Ma, E.S. arXiv:1811.10573

## Regularization and Parallelism for Gauss-Newton



New regularization scheme<sup>1</sup> for Gauss-Newton CP with implicit CG<sup>2</sup>

- Oscillates regularization parameter geometrically between lower and upper thresholds
- Achieves higher convergence likelihood
- More accurate than ALS in applications
- Faster than ALS sequentially and in parallel

<sup>1</sup>Navjot Singh, Linjian Ma, Hongru Yang, and E.S. arXiv:1910.12331

<sup>2</sup>P. Tichavsky, A. H. Phan, and A. Cichocki., 2013



Navjot Singh

## Symmetry in Tensor Contractions



New algebraic transformations for contractions to exploit permutational symmetry<sup>1</sup> and group symmetry<sup>2</sup>, prevalent in quantum chemistry/physics





<sup>1</sup>E.S, J. Demmel, CMAM 2020

<sup>2</sup>collaboration with Yang Gao, Phillip Helms, and Garnet Chan at Caltech

LPNA

### Library for Massively-Parallel Tensor Computations

Cyclops Tensor Framework<sup>1</sup> sparse/dense generalized tensor algebra

- $\bullet\,$  Cyclops is a C++ library that distributes each tensor over MPI
- Used in chemistry (PySCF, QChem)<sup>2</sup>, quantum circuit simulation (IBM/LLNL)<sup>3</sup>, and graph analysis (betweenness centrality)<sup>4</sup>
- Summations and contractions specified via Einstein notation

E["aixbjy"] += X["aixbjy"] - U["abu"]\*V["iju"]\*W["xyu"]

- Best distributed contraction algorithm selected at runtime via models
- Support for Python (numpy.ndarray backend), OpenMP, and GPU
- Simple interface to core ScaLAPACK matrix factorization routines

LPNA

<sup>&</sup>lt;sup>1</sup>https://github.com/cyclops-community/ctf

<sup>&</sup>lt;sup>2</sup>E.S., D. Matthews, J. Hammond, J.F. Stanton, J. Demmel, JPDC 2014

E. Pednault, J.A. Gunnels, G. Nannicini, L. Horesh, T. Magerlein, E. S., E. Draeger, E. Holland, and R. Wisnieff, 2017
E.S., M. Besta, F. Vella, T. Hoefler, SC 2017

# Sparsity in Tensor Contractions



- Cyclops supports sparse representation of tensors<sup>1</sup>
- Choice of representation specified in tensor constructor
- CSR or DCSR<sup>2</sup> (2-index CSF<sup>3</sup>) representation used locally for contractions

<sup>3</sup>S. Smith, G. Karypis 2015

<sup>&</sup>lt;sup>1</sup>E.S., T. Hoefler 2015

<sup>&</sup>lt;sup>2</sup>A. Bulúc, J.R. Gilbert, 2008

# All-at-Once Multi-Tensor Contraction



With sparsity, all-at-once contraction<sup>1</sup> of multiple tensors can be faster<sup>2</sup>.

• Sparse tensor CP decomposition is dominated by MTTKRP

$$u_{ir} = \sum_{j,k} t_{ijk} v_{jr} w_{kr}$$

• All-at-once sparse MTTKRP needs less communication than pairwise

Algorithms for Tensor Computations

• Tensor times tensor product (TTTP) enables sparse residual and CP tensor completion

$$r_{ijk} = \sum_{r} t_{ijk} u_{ir} v_{jr} w_{kr}$$

 Cost and memory footprint reduced asymptotically

<sup>&</sup>lt;sup>1</sup>S. Smith, J. Park, G. Karypis, 2018

<sup>&</sup>lt;sup>2</sup>Zecheng Zhang, Xiaoxiao Wu, Naijing Zhang, Siyuan Zhang, and E.S. arXiv:1910.02371

## Conclusion

- Communication avoiding matrix factorizations
  - ${\, \bullet \, }$  3D algorithms move  ${\cal O}(p^{1/6})$  less data on p processors
  - Faster algorithms for TRSM, Cholesky, LU, QR, and symmetric eigensolve
- Optimization methods for tensor decomposition
  - Pairwise perturbation approximates ALS sweeps with asymptotically less cost
  - Gauss-Newton methods with new regularization scheme more effective than ALS
- Algorithms and software for sparse/symmetric tensors
  - Algebraic reorganizations to exploit permutational and group symmetry
  - Cyclops library for dense/sparse generalized distributed tensor algebra in C++/Python

## Acknowledgements

- Laboratory for Parallel Numerical Algorithms (LPNA) at University of Illinois
- Key external collaborators: James Demmel, Torsten Hoefler, Garnet Chan
- NSF OAC CSSI funding (award #1931258)
- DOE Computational Science Graduate Fellowship
- Stampede2 resources at TACC via XSEDE



L P. N A @CS@Illinois



http://lpna.cs.illinois.edu