

Scaling Numerical Algorithms and Software via Improved Performance Modeling

Edgar Solomonik

LPNA @CS@Illinois

Department of Computer Science
University of Illinois at Urbana-Champaign

ICL Lunch Talk
University of Tennessee, Knoxville

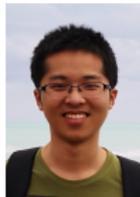
Laboratory for Parallel Numerical Algorithms

Focus today

- communication avoiding algorithms for dense linear algebra
- Cyclops library for tensor contractions
- inexact autotuning via critical path profiling
- performance modeling via tensor completion

See <http://lpna.cs.illinois.edu> for our group's other work

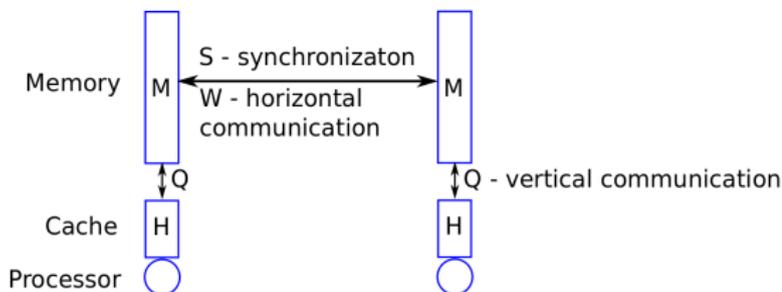
- quantum simulation
- linear system solvers for interior point methods
- optimization and sketching for tensor decompositions/networks
- parallel sorting



Beyond computational complexity

Algorithms should minimize communication, not just computation

- communication and synchronization cost more energy than flops
- two types of communication (data movement):



- **vertical** (intranode memory–cache)
- **horizontal** (internode network transfers)
- parallel algorithm design involves tradeoffs: computation vs communication vs synchronization
- parameterized algorithms provide optimality and flexibility

Cost model for parallel algorithms

We use the **Bulk Synchronous Parallel (BSP) model** (L.G. Valiant 1990)

- execution is subdivided into S supersteps, each associated with a global **synchronization** (cost α)
- at the start of each superstep, processors interchange messages, then they perform local computation
- if the **maximum amount of data** sent or received by any process is w_i (work done is f_i and amount of memory traffic is q_i) at superstep i then the BSP time is

$$T = \sum_{i=1}^S \alpha + w_i \cdot \beta + q_i \cdot \nu + f_i \cdot \gamma = O(S \cdot \alpha + W \cdot \beta + Q \cdot \nu + F \cdot \gamma)$$

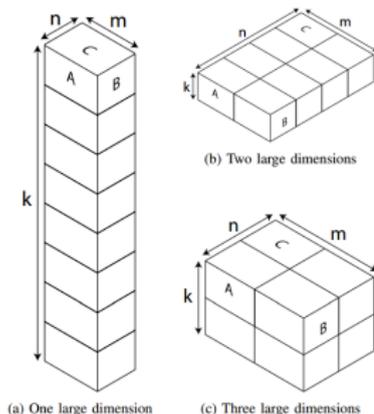
where typically $\alpha \gg \beta \gg \nu \gg \gamma$

- we mention vertical communication cost only when it exceeds $Q = O(F/\sqrt{H} + W)$ where H is cache size

Communication complexity of matrix multiplication

Multiplication of $\mathbf{A} \in \mathbb{R}^{m \times k}$ and $\mathbf{B} \in \mathbb{R}^{k \times n}$ can be done in $O(1)$ supersteps with **communication cost** $W = O\left(\left(\frac{mnk}{p}\right)^{2/3}\right)$ provided sufficient memory and sufficiently large p

- when $m = n = k$, 3D blocking gets $O(p^{1/6})$ improvement over 2D¹
- when m, n, k are unequal, need appropriate processor grid²



¹J. Berntsen, Par. Comp., 1989; A. Aggarwal, A. Chandra, M. Snir, TCS, 1990; R.C. Agarwal, S.M. Balle, F.G. Gustavson, M. Joshi, P. Palkar, IBM, 1995; F.W. McColl, A. Tiskin, Algorithmica, 1999; ...

²J. Demmel, D. Eliahu, A. Fox, S. Kamil, B. Lipshitz, O. Schwartz, O. Spillinger 2013

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**.

triangular solve	square TRSM \checkmark^1	rectangular TRSM \checkmark^2
LU with pivoting	pairwise pivoting \checkmark^3	tournament pivoting \checkmark^4
QR factorization	Givens on square \checkmark^3	Householder on rect. \checkmark^5
sym. eig.	eigenvalues only \checkmark^5	eigenvectors \times

\checkmark means costs attained (synchronization within polylog factors).

¹B. Lipshitz, MS thesis 2013

²T. Wicky, E.S., T. Hoefler, IPDPS 2017

³A. Tiskin, FGCS 2007

⁴E.S., J. Demmel, EuroPar 2011

⁵E.S., G. Ballard, T. Hoefler, J. Demmel, SPAA 2017

Characteristics of 3D Algorithms

- recursive formulations possible for many problems
 - Tiskin's algorithms for Cholesky, LU with pairwise pivoting, QR
 - some are complicated, latter two recurse on slanted 2:1 panels
- two-level (logical) blocking
 - yields 3D algorithms for LU, QR and reduces vertical communication
 - also used in full-to-band reduction for symmetric eigensolve
- successive band reduction
 - classically used to reduce vertical comm.¹
 - $\log p$ reduction stages needed to obtain 3D symmetric eigensolve²
- alternative numerical formulations
 - different pivoting (pairwise/tournament) for LU, polar decomposition³
 - Cholesky-QR2 and variants^{4,5}, triangular inversion (log depth)⁶

¹C.H. Bischof, B. Lang, X. Sun, ACM TOMS'00

²E.S., G. Ballard, J. Demmel, T. Hoefler, SPAA'17

³Y. Nakatsukasa, N.J. Higham, SISC'13. H. Ltaief, D. Sukkari, A. Esposito, Y. Nakatsukasa, D. Keyes, ACM TOPC'19.

⁴T. Fukaya, R. Kannan, Y. Nakatsukasa, Y. Yamamoto, Y. Yanagisawa, SISC'20

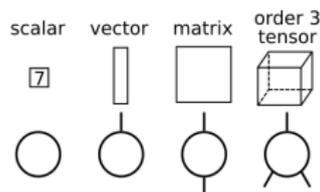
⁵E. Hutter, E.S., IPDPS'19

⁶T. Wicky, T. Hoefler, E.S., IPDPS'17

Tensors

A **tensor** is a collection of elements

- its **dimensions** define the size of the collection
- its **order** is the number of different dimensions
- specifying an index along each tensor **mode** defines an element of the tensor



A few examples of tensors are

- Order 0 tensors are scalars, e.g., $s \in \mathbb{R}$
- Order 1 tensors are vectors, e.g., $v \in \mathbb{R}^n$
- Order 2 tensors are matrices, e.g., $A \in \mathbb{R}^{m \times n}$
- An order 3 tensor with dimensions $s_1 \times s_2 \times s_3$ is denoted as $\mathcal{T} \in \mathbb{R}^{s_1 \times s_2 \times s_3}$ with elements t_{ijk} for $i \in \{1, \dots, s_1\}, j \in \{1, \dots, s_2\}, k \in \{1, \dots, s_3\}$

Tensor Contractions

A **tensor contraction** describes a set of products and sums of elements from two tensors

tensor contraction	formula
inner product	$w = \sum_i u_i v_i$
outer product	$w_{ij} = u_i v_{ij}$
pointwise product	$w_i = u_i v_i$
Hadamard product	$w_{ij} = u_{ij} v_{ij}$
matrix multiplication	$w_{ij} = \sum_k u_{ik} v_{kj}$
batched mat.-mul.	$w_{ijl} = \sum_k u_{ikl} v_{kjl}$
tensor times matrix	$w_{ilk} = \sum_j u_{ijk} v_{lj}$

Tensor contractions are prevalent in quantum chemistry methods

Library for Massively-Parallel Tensor Contractions

Cyclops Tensor Framework¹: sparse/dense generalized tensor algebra

- Cyclops is a C++ library that distributes each tensor over MPI
- Used in chemistry (PySCF, QChem, CC4S)², quantum circuit simulation (by IBM/LLNL)³, and graph analysis (betweenness centrality⁴, minimum spanning tree⁵)

- 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 CUDA

¹<https://github.com/cyclops-community/ctf>

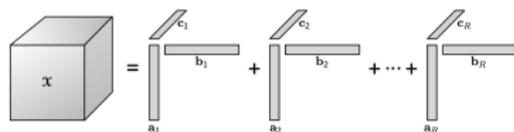
²E.S., D. Matthews, J. Hammond, J.F. Stanton, J. Demmel, JPDC 2014

³E. Pednault, J.A. Gunnels, G. Nannicini, L. Horesh, T. Magerlein, E.S., E. Draeger, E. Holland, and R. Wisnieff, 2017

⁴E.S., M. Besta, F. Vella, T. Hoefer, SC 2017

⁵T. Baer, R. Kanakagiri, E.S., SIAM PP 2022

CP Tensor Decomposition



- For a tensor $\mathcal{T} \in \mathbb{R}^{n \times n \times n}$, the CP decomposition^{1,2} is defined by matrices A , B , and C such that

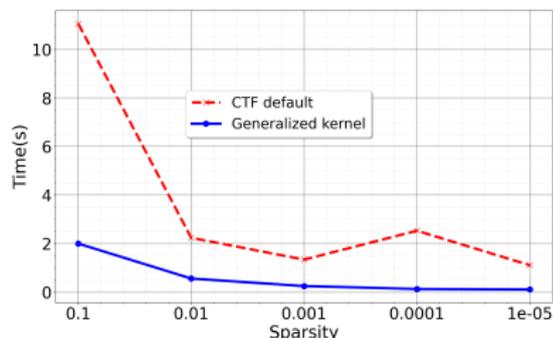
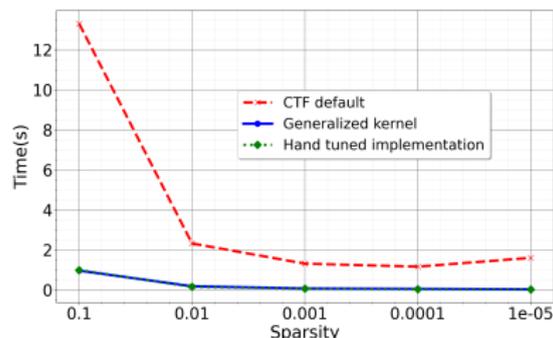
$$t_{ijk} = \sum_{r=1}^R a_{ir} b_{jr} c_{kr}$$

- low-rank CP decomposition is widely used for compression and multi-way data analysis
- high-rank CP decomposition is useful in quantum simulation, search for fast bilinear algorithms

¹F.L. Hitchcock, Studies in Applied Mathematics 1927

²T. Kolda and B. Bader, SIAM Review 2009

Recent and Ongoing Cyclops Developments



- All-at-once contraction for sparse tensor times dense tensor network
 - Driven by tensor completion¹ and quasi-robust density fitting²
 - Generalizes MTTKRP (common kernel for CP), TTMc (for Tucker) and other kernels arising in sparse tensor decomposition and completion
 - Working on integration with linear (least-squares) solves
 - Select best loop-nest / integrate BLAS based on performance model

¹N. Singh, Z. Zhang, X. Wu, N. Zhang, S. Zhang, and E.S., JPDC'22

²D.P. Tew, The Journal of Chemical Physics 2018

Accelerating Automatic Tuning

Autotuning – searching for fastest program variant by benchmarking

- widely used for performance tuning today
- computationally expensive, especially if different variants wanted for different inputs

Performance of parallel programs is hard to predict

- communication bottlenecks dependent on architecture
- idle time and load imbalance arise due to dependencies

But many programs largely consist of easy-to-model subkernels

- parallel QR codes repeatedly execute local MMs of similar sizes

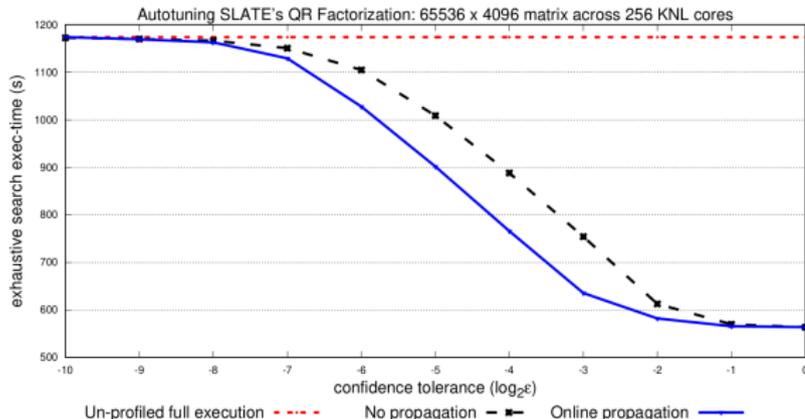
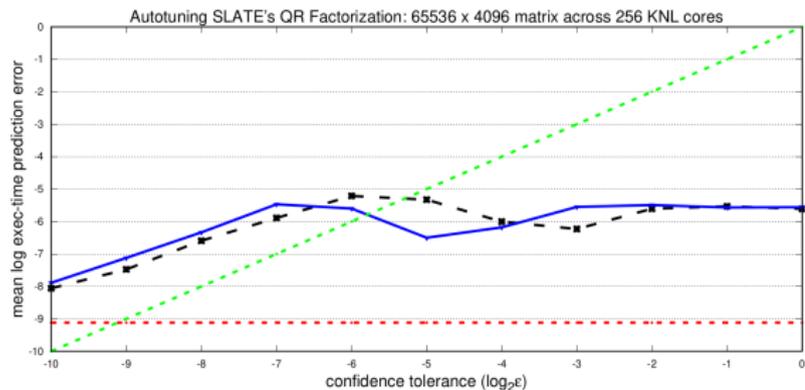
Approximate Automatic Tuning

Critter¹ leverages critical path profiling to efficiently perform approximate autotuning in dense linear algebra

- designed to track critical path communication/synchronization costs
- approximates critical path execution time by predicting subkernel performance
- the more times a subkernel is executed, the more statistical accuracy in overall prediction
- propagates local kernel timings to other processors to reduce cost
- automatically profiles MPI routines and identifies basic BLAS routines

¹E. Hutter, E.S. IPDPS'21

Approximate Automatic Tuning for QR



Performance Modeling

- selection of best program variants often guided by performance modeling
 - for autotuning, performance models can guide/accelerate
- in practice, performance models may be
 - semi-analytic or derived from program structure
 - learned from sample of benchmark timings
- for complex programs, input and parameter space is high-dimensional
 - multi-task learning problem
 - runtime data is noisy and partially complete

Prediction Model and Accuracy

- Given d input/tuning parameters and executions times $T(p_{i_1}, \dots, p_{i_d}) \in \mathbb{R}_+$ for $(i_1, \dots, i_d) \in \Omega$, seek concise model T for all other parameters
- For example, seek to model $m \times k$ by $k \times n$ matrix multiplication
 - analytically, a simple model may be

$$T_{\text{MM}}(m, n, k) = \gamma \cdot mnk + \beta \cdot (mn + nk + mk) + \alpha$$

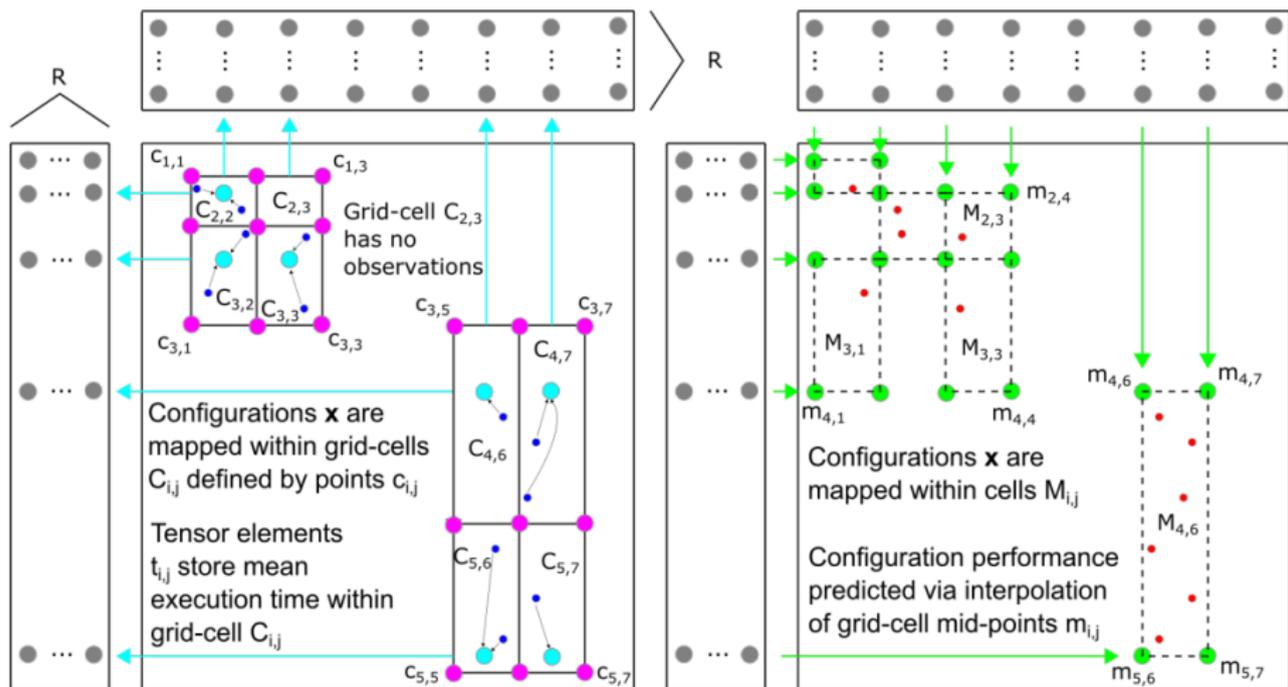
with α , β , and γ tuned

- Generally, we want the prediction to be accurate in scale, namely we want to minimize error in $\log(T)$
- We apply CP tensor completion for modeling, e.g.,

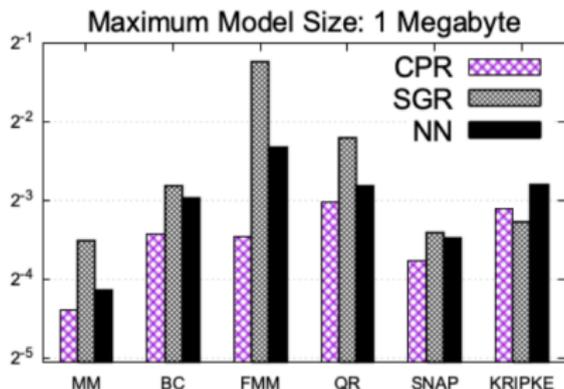
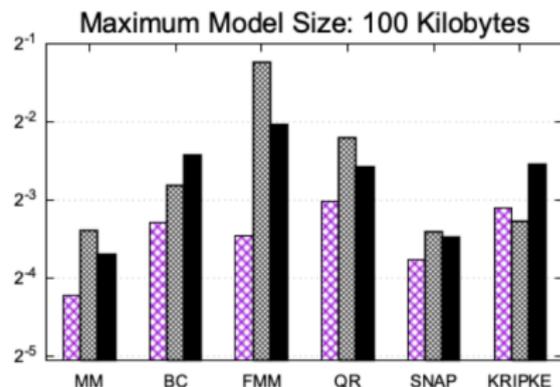
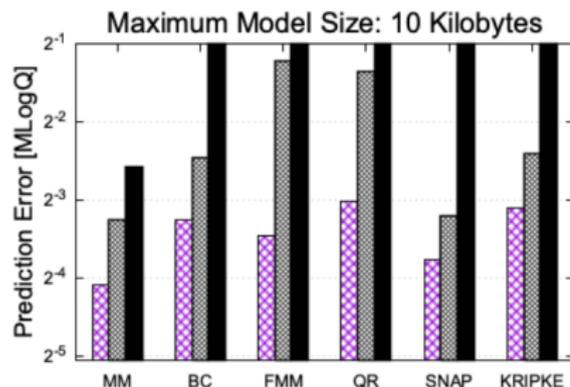
$$\min_{\mathbf{x}=\text{CP}(A,B,C)} \sum_{(u,v,w) \in \Omega} (x_{uvw} - \log(T_{\text{MM}}(m_u, n_v, k_w)))^2 + \lambda \dots$$

- We find the log significantly improves effective prediction accuracy, and is more efficient than using a log loss function

Discretization and Interpolation



Discretization and Interpolation



Conclusion and Acknowledgements

- For details on tensor completion modeling, see <https://solomonik.cs.illinois.edu/preprints/HS22.pdf>
- NSF awards #1839204 (RAISE-TAQS), #1931258 (CSSI), #1942995 (CAREER), the DOE CSGF program, and the DOE MMICC program
- Stampede2 resources at TACC via XSEDE
- See <https://lpna.cs.illinois.edu> for all further info



LPNA @ CS @ Illinois

