


# High Performance Tensor Computations

Edgar Solomonik

 @CS@Illinois

Department of Computer Science  
University of Illinois at Urbana-Champaign

SPCL\_Bcast(COMM\_WORLD) seminar

# Laboratory for Parallel Numerical Algorithms

Recent/ongoing research topics  
(\*covered today)

- parallel matrix computations
  - matrix factorizations
  - eigenvalue problems
  - preconditioners
- 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\*



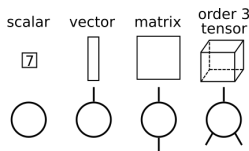
LPNA @ CS@Illinois



<http://lpna.cs.illinois.edu>

# Tensor Contractions

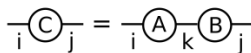
- A tensor of order  $N$  has  $N$  modes and dimensions  $s \times \cdots \times s$



- Two or more tensors can be contracted together in various ways

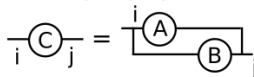
Matrix Multiplication

$$c_{ij} = \sum_k a_{ik} b_{kj}$$



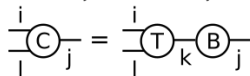
Hadamard Product

$$c_{ij} = a_{ij} b_{ij}$$



Tensor Times Matrix (TTM)

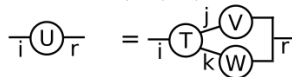
$$c_{ijl} = \sum_k t_{ikl} b_{kj}$$



MTTKRP

(matricized tensor times Khatri-Rao product)

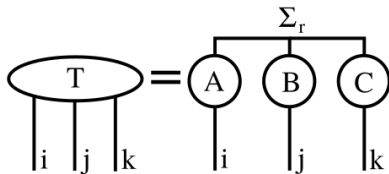
$$u_{ir} = \sum_{jk} t_{ijk} v_{jr} w_{kr}$$



# Tensor Decompositions

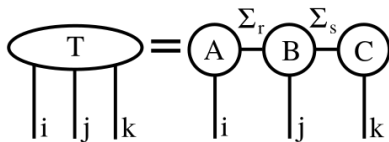
- Canonical polyadic (CP) tensor decomposition<sup>1</sup>

$$t_{ijk} = \sum_{r=1}^R u_{ir} v_{jr} w_{kr}$$

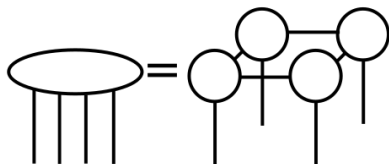


- 1D tensor network / Matrix product state (MPS) / tensor train (TT) decomposition

$$t_{ijk} = \sum_r \sum_s u_{ir} v_{rj} w_{sk}$$



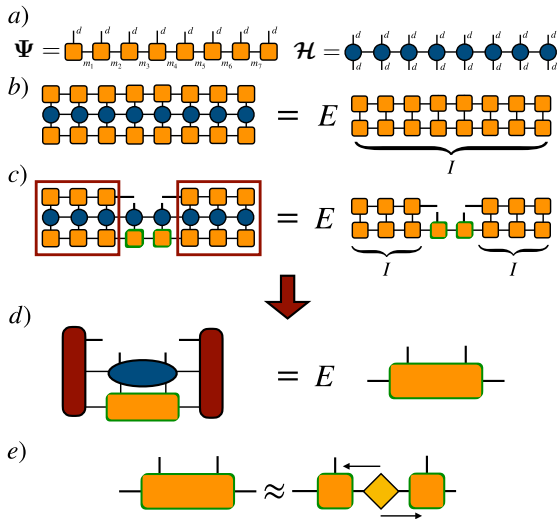
- 2D tensor network / projected entangled pair state (PEPS)



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



# Tensor Network Methods



# Applications

- Tensor Decompositions

- data mining e.g., high-order clustering, typically low-rank decomposition augmented with constraints
- model/data compression e.g., neural networks and quantum chemistry, relatively high rank needed for accuracy
- discovery of bilinear algorithms (e.g., Strassen's algorithm for matrix multiplication), small sparse tensors with relatively high rank

# Applications

- Tensor Decompositions

- data mining e.g., high-order clustering, typically low-rank decomposition augmented with constraints
- model/data compression e.g., neural networks and quantum chemistry, relatively high rank needed for accuracy
- discovery of bilinear algorithms (e.g., Strassen's algorithm for matrix multiplication), small sparse tensors with relatively high rank

- Tensor Networks

- eigenvalue and least squares problems where system of equations and vector are represented by low rank tensor networks
- prevalent in simulation of quantum systems (spins, electrons, qubits) given Hamiltonian or quantum circuit description

# Applications

- Tensor Decompositions

- data mining e.g., high-order clustering, typically low-rank decomposition augmented with constraints
- model/data compression e.g., neural networks and quantum chemistry, relatively high rank needed for accuracy
- discovery of bilinear algorithms (e.g., Strassen's algorithm for matrix multiplication), small sparse tensors with relatively high rank

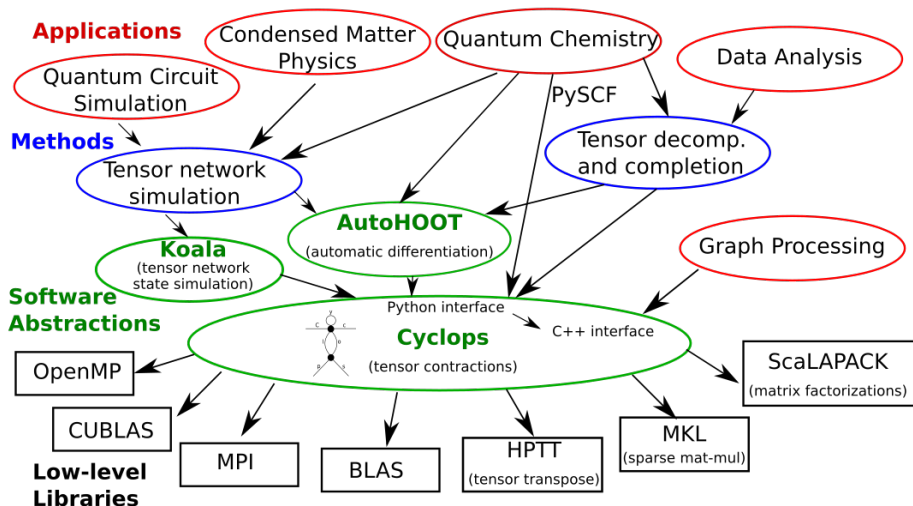
- Tensor Networks

- eigenvalue and least squares problems where system of equations and vector are represented by low rank tensor networks
- prevalent in simulation of quantum systems (spins, electrons, qubits) given Hamiltonian or quantum circuit description

- Tensor Contractions

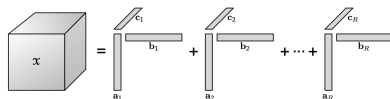
- prevalent within tensor network and tensor decomposition methods
- also arise in high-accuracy quantum chemistry methods (e.g., coupled cluster), where tensors often have symmetries

# Software Abstractions for Tensor Computations



# CP Tensor Decomposition Algorithms

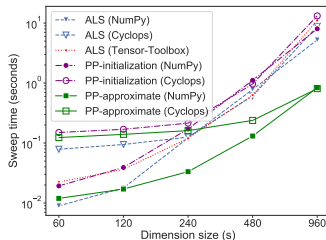
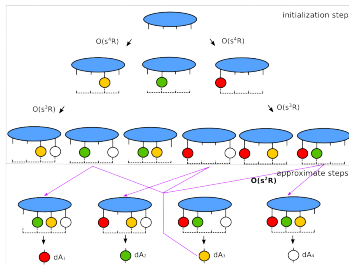
- Tensor of **order**  $N$  has  $N$  **modes** and **dimensions**  $s \times \cdots \times s$
- Canonical polyadic (**CP**) tensor decomposition<sup>1</sup>



- Alternating least squares (**ALS**) is most widely used method
  - Optimize one factor matrix at a time, yielding quadratic optimization subproblems
  - Achieves monotonic linear convergence
- **Gauss-Newton** method is an emerging alternative
  - Optimizes all factor matrices at once by quadratic approximation of nonlinear objective function
  - Non-monotonic, but can achieve quadratic convergence

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

# Pairwise Perturbation Algorithm



New algorithm: **pairwise perturbation (PP)**<sup>1</sup> approximates ALS

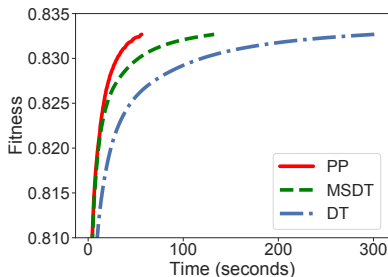
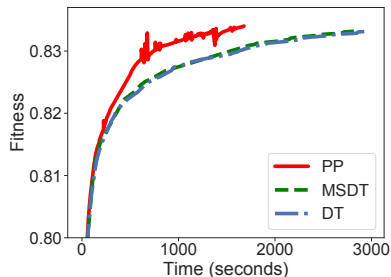
- based on perturbative expansion of ALS update to approximate MTTKRP
- approximation is accurate when ALS updates stagnate
- rank  $R < s^{N-1}$  CP decomposition:
  - ALS sweep cost  $O(s^N R) \Rightarrow O(s^2 R)$ , up to 33x speed-up



Linjian Ma

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

# Parallel Pairwise Perturbation Algorithm



Effective parallelization by decomposing MTTKRP into local MTTKRPs <sup>1</sup>

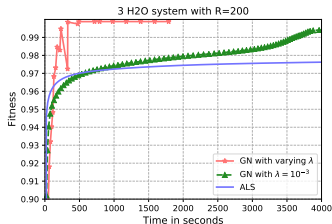
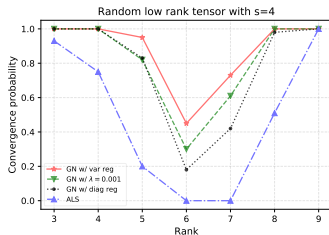
$$U = \text{MTTKRP}(\mathcal{T}, V, W) \Rightarrow U_i = \sum_{j,k} \text{MTTKRP}(\mathcal{T}_{ijk}, V_j, W_k)$$

- processor  $(i, j, k)$  owns  $\mathcal{T}_{ijk}$ ,  $V_j$ , and  $W_k$
- pairwise perturbation can be used to approximate local MTTKRPs, reducing communication cost
- multi-sweep dimension-tree (MSDT) amortizes terms across sweeps

<sup>1</sup>L. Ma, E.S. to appear on arXiv, October 2020.



# 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



Navjot Singh

<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

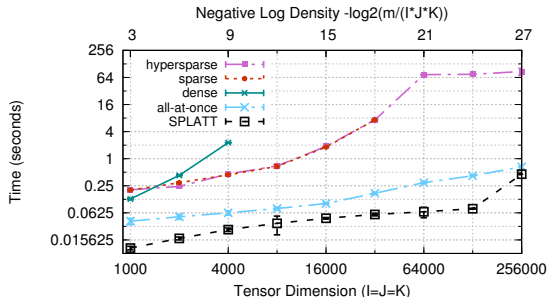
# Sparse Tensor Decomposition

- Sparse tensor decomposition is dominated by MTTKRP

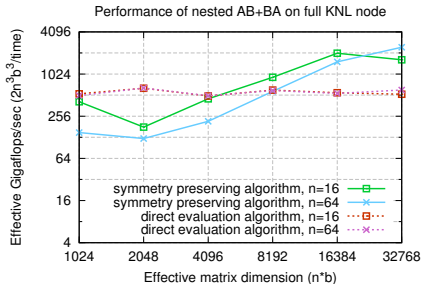
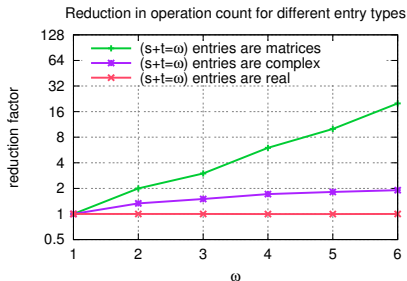
$$u_{ir} = \sum_{j,k} t_{ijk} v_{jr} w_{kr}$$

- Sparse MTTKRP can be done faster all-at-once than by contracting two tensors at a time

Cyclops MTTKRP with  $m=1B$ ,  $R=50$  on 4096 Cores of Stampede2



# Permutational Symmetry in Tensor Contractions



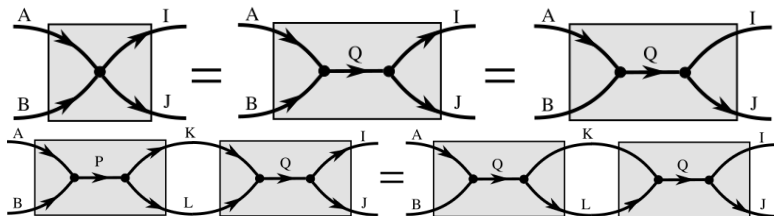
New contraction algorithms reduce cost via permutational symmetry<sup>1</sup>

- Symmetry is hard to use in contraction e.g.  $y = Ax$  with  $A$  symmetric
- For contraction of order  $s + v$  and  $v + t$  tensors to produce an order  $s + t$  tensor, previously known approaches reduce cost by  $s!t!v!$
- New algorithm reduces number of *products* by  $\omega!$  where  $\omega = s + t + v$ , leads to same reduction in *cost* for partially-symmetric contractions

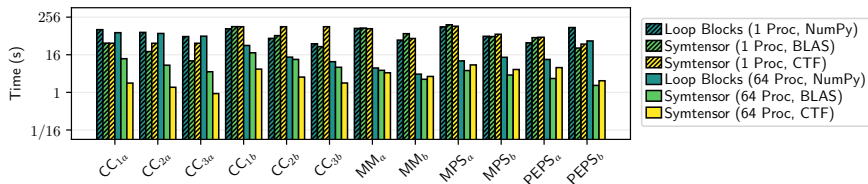
$$C = AB + BA \Rightarrow c_{ij} = \sum_k [(a_{ij} + a_{ik} + a_{jk}) \cdot (b_{ij} + b_{ik} + b_{jk})] - \dots$$

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

# Group Symmetry in Tensor Contractions



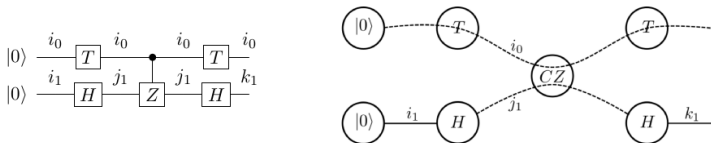
New contraction algorithm, *irreducible representation alignment* uses new reduced form to handle group symmetry (momentum conservation, spin, quantum numbers, etc.) without looping over blocks or sparsity<sup>1</sup>



<sup>1</sup>Y. Gao, P. Helms, G. Chan, and E.S., arXiv:2007.08056

# Quantum Circuit Simulation with Tensor Networks

- A quantum circuit is a direct description of a tensor network<sup>1</sup>



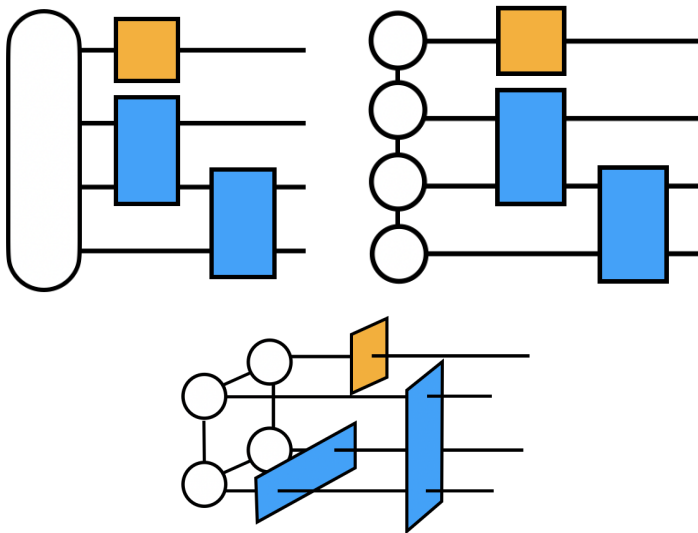
- Why use HPC to (approximately) simulate quantum circuits?
  - enable development/testing/tuning of larger quantum circuits
  - understand approximability of different quantum algorithms
  - quantify sensitivity of algorithms to noise/error
  - potentially enable new hybrid quantum-classical algorithms
- Cyclops utilized to simulate 49-qubit circuits by IBM+LLNL team via direct contraction<sup>2</sup> and by another team from via exact PEPS evolution/contraction<sup>3</sup>

<sup>1</sup>Markov and Shi SIAM JC 2007

<sup>2</sup>Pednault et al. arXiv:1710.05867

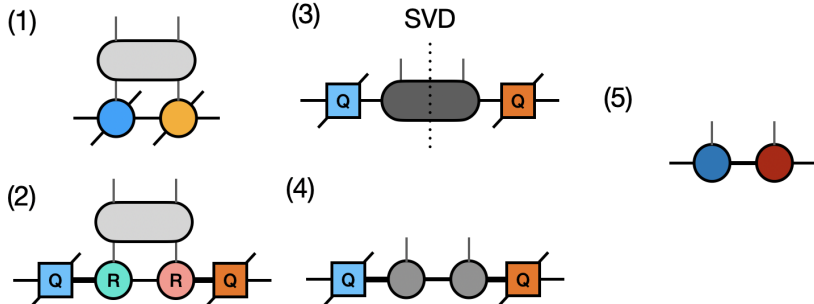
<sup>3</sup>Guo et al. Phys Rev Letters, 2019

# Tensor Network State Simulation



# Approximate Application of Two-Site Operators

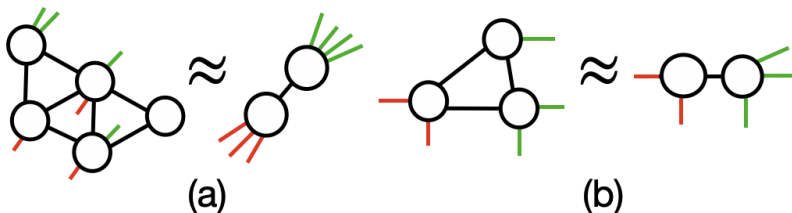
- Consider application of a two-site operator on neighboring PEPS sites
- Simple update (QR-SVD)* algorithm:



- We provide an efficient distributed implementation of QR-SVD
- This operation is an instance of what we'll refer to as `einsumsvd` and QR-SVD is one algorithm/implementation

# Implicit Randomized einsumsvd

- The einsumsvd primitive will also enable effective algorithms for PEPS contraction

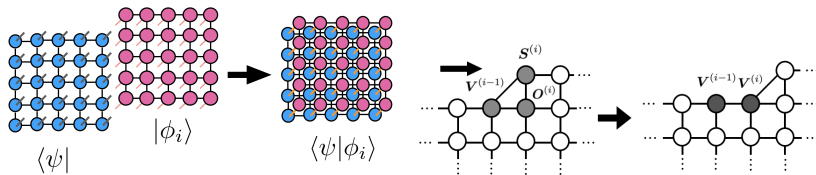


- An efficient general implementation is to leverage randomized SVD / orthogonal iteration, which iteratively computes a low-rank SVD by a matrix–matrix product that can be done implicitly via tensor contractions



# PEPS Contraction

- Exact contraction of PEPS is  $\#P$ -complete, so known methods have exponential cost in the number of sites
- PEPS contraction is needed to compute expectation values
- Boundary contraction* is common for finite PEPS and can be simplified with einsumsvd

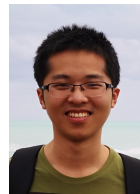


- We introduce a new library, Koala<sup>1</sup>, for high-performance simulation of quantum circuits and time evolution with PEPS<sup>2</sup>

```

1  from koala import peps, Observable
2  from tensorbackends.interface import ImplicitRandomizedSVD
3
4  # Create a 2-by-3 PEPS in distributed memory using CTF
5  qstate = peps.computational_zeros(nrow=2, ncol=3, backend='ctf')
6
7  # Construct operators and apply them to the quantum state
8  Y = qstate.backend.astensor([0,-1j,1j,0]).reshape(2,2)
9  CX = qstate.backend.astensor([1,0,0,0,0,1,0,0,0,0,1,0,0,1,0,])
10 CX = CX.reshape(2,2,2,2)
11
12 qstate.apply_operator(Y, [1])
13 qstate.apply_operator(CX, [1,4], update_option=peps.QRUpdate(rank=2))
14
15 # Construct an observable and calculate the expectation with IBMPs
16 observable = Observable.ZZ(3, 4) + 0.2 * Observable.X(1)
17 result = qstate.expectation(
18     observable, use_cache=True,
19     contract_option=peps.BMPS(ImplicitRandomizedSVD(rank=4)),
20 )

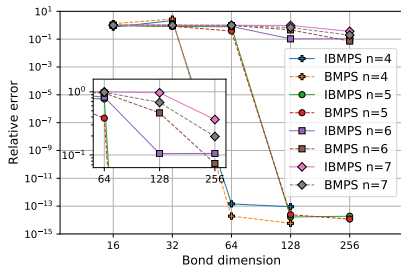
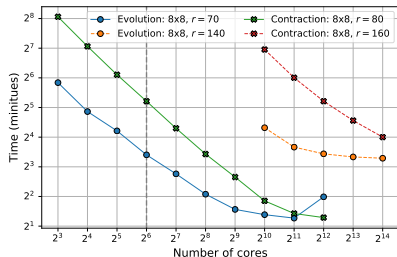
```



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

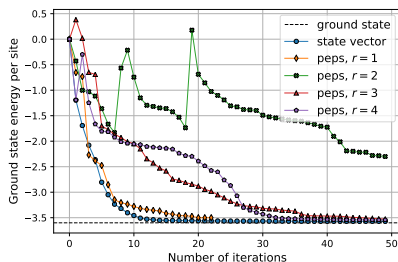
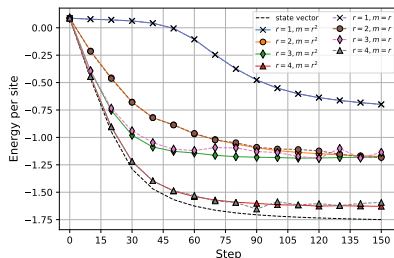
<sup>2</sup>Yuchen Pang, Tianyi Hao, Annika Dugad, Yiqing Zhou, and E.S., to appear in proceedings of SC 2020, arXiv:2006.15234.

# PEPS Benchmark Performance



- Koala achieves good parallel scalability for approximate gate application (evolution) and contraction
- Approximation can be effective even for adversarially-designed circuits such as Google's random quantum circuit model (figure on right)

# PEPS Accuracy for Quantum Simulation



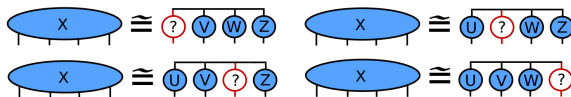
- ITE code achieves improvable accuracy with increased PEPS bond dimension, but approximation in PEPS contraction is not variational
- Variational quantum eigensolver (VQE), which represents a wavefunction using a parameterized circuit  $U(\theta)$  and minimizes

$$\langle U(\theta) | H | U(\theta) \rangle ,$$

also achieves improvable accuracy with higher PEPS bond dimension

# Automatic Differentiation for Tensor Computations

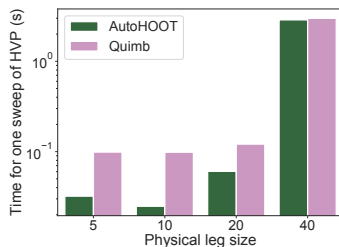
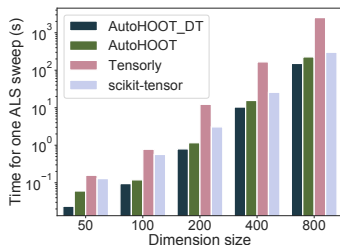
- Tensor network and tensor decomposition methods all typically based on applying Newton's method on a sequence of subsets of variables



- Automatic differentiation (AD) in principle enables automatic generation of these methods
- However, existing AD tools such as Jax (used by TensorFlow) are designed for deep learning and are ineffective for more complex tensor computations
  - these focus purely on first order optimization via Jacobian-vector products
  - unable to propagate tensor algebra identities such as  $(A \otimes B)^{-1} = A^{-1} \otimes B^{-1}$  to generate efficient code

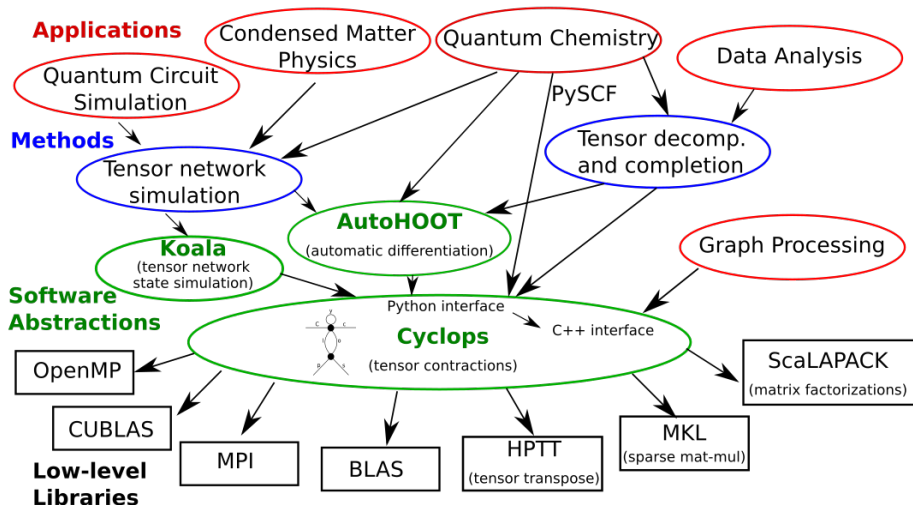
# AutoHOOT: Automatic High-Order Optimization for Tensors

- AutoHOOT<sup>1</sup> provides a tensor-algebra centric AD engine
- Designed for einsum expressions and alternating minimization common in tensor decomposition and tensor network methods
- Python-level AD is coupled with optimization of contraction order and caching of intermediates
- Generates code for CPU/GPU/supercomputers using high-level back-end interface to tensor contractions



<sup>1</sup>Linjian Ma, Jiayuan Ye, and E.S. arXiv:2005.04540, 2020

# Software Abstractions for Tensor Computations



# Acknowledgements

- Laboratory for Parallel Numerical Algorithms (LPNA) at University of Illinois, [lpna.cs.illinois.edu](http://lpna.cs.illinois.edu) and collaborators
- Funding from NSF awards: #1839204 (RAISE-TAQS), #1931258 (CSSI), #1942995 (CAREER)
- Stampede2 resources at TACC via XSEDE



L · P · N A @ CS @ Illinois



<http://lpna.cs.illinois.edu>



# Backup slides

# Library for Massively-Parallel Tensor Computations

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

- 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

---

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

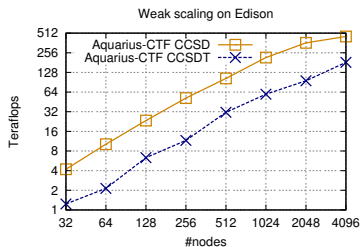
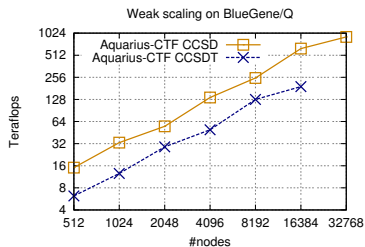
<sup>2</sup>E.S., D. Matthews, J. Hammond, J.F. Stanton, J. Demmel, JPDC 2014

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

<sup>4</sup>E.S., M. Besta, F. Vella, T. Hoefer, SC 2017

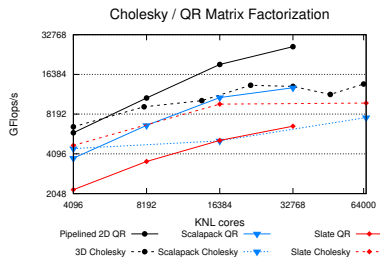
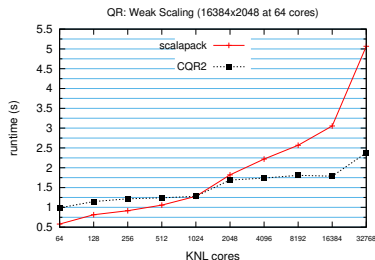
# Electronic structure calculations with Cyclops

CCSD up to 55 (50) water molecules with cc-pVDZ  
CCSDT up to 10 water molecules with cc-pVDZ



compares well to NWChem (up to 10x speed-ups for CCSDT)

# Faster Parallel Algorithms for Cholesky and QR



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

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

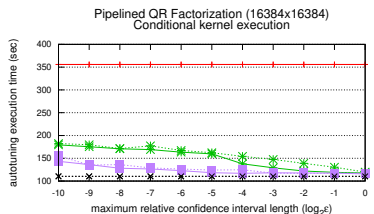
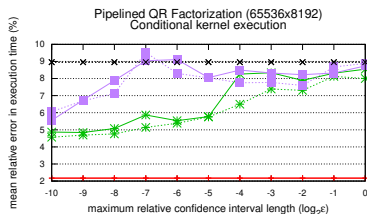


Edward Hutter

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

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

# Accelerating Autotuning using Critter



New tool **critter** accelerates autotuning by conditional execution<sup>1</sup>

- Avoids execution of computation and communication kernels if their performance is predictable using past observations
- Leverages critical path profiling to determine needed model accuracy for each kernel



Edward Hutter

<sup>1</sup>E. Hutter, Edgar Solomonik, to appear on arXiv