

Scalable Tensor Algorithms for Scientific Computing

Edgar Solomonik

Department of Computer Science
University of Illinois at Urbana-Champaign

PMAA 2018

June 28, 2018

 @CS@Illinois

A library for parallel tensor computations

Cyclops Tensor Framework (CTF)¹

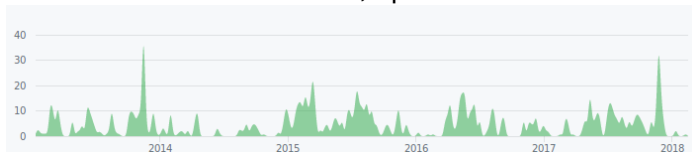
- distributed-memory symmetric/sparse/dense tensor objects

```
Matrix<int> A(n, n, AS|SP, World(MPI_COMM_WORLD));  
Tensor<float> T(order, is_sparse, dims, syms, ring, world);  
T.read(...); T.write(...); T.slice(...); T.permute(...);
```

- parallel contraction/summation of tensors

```
Z["abij"] += V["ijab"];  
B["ai"]    = A["aiai"];  
T["abij"]  = T["abij"]*D["abij"];  
W["mnij"] += 0.5*W["mnef"]*T["efij"];  
Z["abij"] -= R["mnje"]*T3["abeimn"];  
M["ij"] += Function<>([](double x){ return 1/x; })(v["j"]);
```

- ~2000 commits since 2011, open source since 2013



¹E.S., D. Matthews, J.R. Hammond, J. Demmel, JPDC 2014

Electronic structure calculations with cyclops

Extracted from [Aquarius](#) (lead by [Devin Matthews](#))

<https://github.com/devinamatthews/aquarius>

```
FMI["mi"]      += 0.5*WMNEF["mnef"]*T2["efin"];
WMNIJ["mnij"]  += 0.5*WMNEF["mnef"]*T2["efij"];
FAE["ae"]      -= 0.5*WMNEF["mnef"]*T2["afmn"];
WAMEI["amei"]  -= 0.5*WMNEF["mnef"]*T2["afin"];

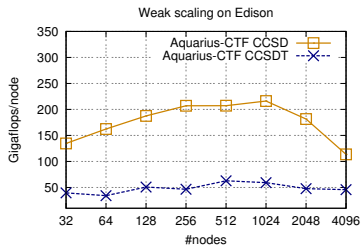
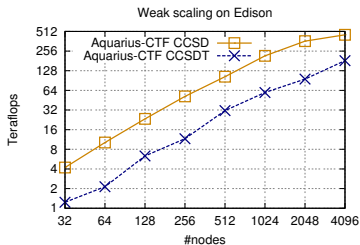
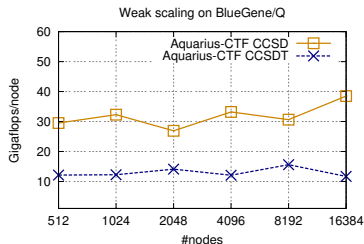
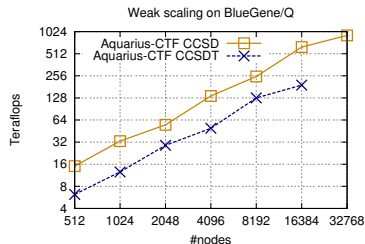
Z2["abij"]     = WMNEF["ijab"];
Z2["abij"]     += FAE["af"]*T2["fbij"];
Z2["abij"]     -= FMI["ni"]*T2["abnj"];
Z2["abij"]     += 0.5*WABEF["abef"]*T2["efij"];
Z2["abij"]     += 0.5*WMNIJ["mnij"]*T2["abmn"];
Z2["abij"]     -= WAMEI["amei"]*T2["ebmj"];
```

- CTF has been integrated with **QChem**, **VASP**, and **Psi4**
- Is also being used for other applications, e.g. by IBM+LLNL collaboration to perform 49-qubit quantum circuit simulation

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)

```
Tensor<> Ea, Ei, Fab, Fij, Vabij, Vijab, Vabcd, Vijkl, Vaibj;  
... // compute above 1-e and 2-e integrals
```

```
Tensor<> T(4, Vabij.lens, *Vabij.world);  
T["abij"] = Vabij["abij"];
```

```
divide_EaEi(Ea, Ei, T);
```

```
Tensor<> Z(4, Vabij.lens, *Vabij.world);  
Z["abij"] = Vijab["ijab"];  
Z["abij"] += Fab["af"]*T["fbij"];  
Z["abij"] -= Fij["ni"]*T["abnj"];  
Z["abij"] += 0.5*Vabcd["abef"]*T["efij"];  
Z["abij"] += 0.5*Vijkl["mnij"]*T["abmn"];  
Z["abij"] += Vaibj["amei"]*T["ebmj"];
```

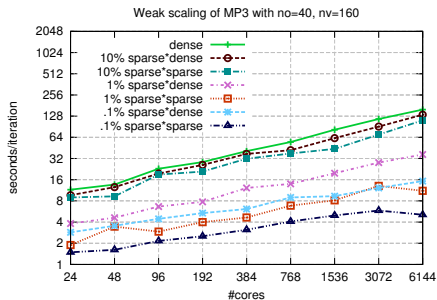
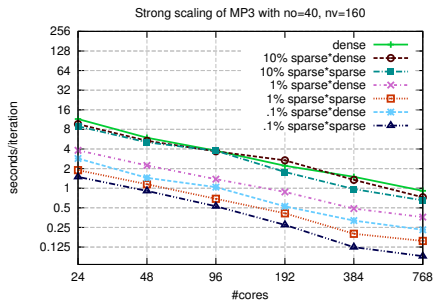
```
divide_EaEi(Ea, Ei, Z);
```

```
double MP3_energy = Z["abij"]*Vabij["abij"];
```

Sparse MP3 code

Strong and weak scaling of sparse MP3 code, with

(1) dense V and T (2) sparse V and dense T (3) sparse V and T



Custom tensor element types

Cyclops permits arbitrary element types and custom functions

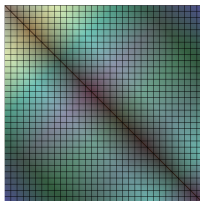
- CombBLAS/GraphBLAS-like functionality
- See examples for SSSP, APSP, betweenness centrality, MIS, MIS-2
- Functionality to handle serialization of pointers within user-defined types is under development
- Can already do block-sparsity via sparse tensor (local) of dense tensors (parallel)

```
Matrix< Matrix<> > C(nblk, nblk, SP, self_world, tmon);
```

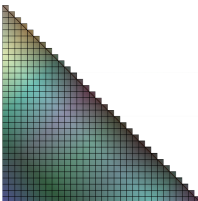
```
C["ij"] = Function< Matrix<> >(  
    [](Matrix<> mA, Matrix<> mB){  
        mC["ij"] += mA["ik"]*mB["kj"];  
        return mC;  
    })  
(A["ik"], B["kj"]);
```

Symmetry and sparsity by cyclicity

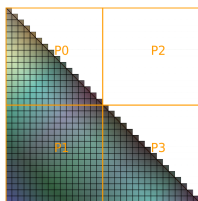
Symmetric matrix



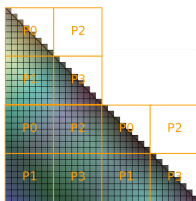
Unique part of symmetric matrix



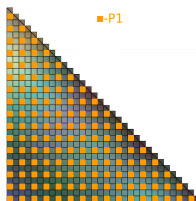
Naive blocked layout



Block-cyclic layout

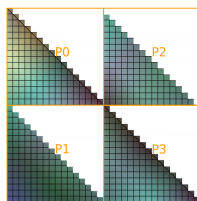


Cyclic layout



~

Improved blocked layout



for sparse tensors, a cyclic layout provides a load-balanced distribution

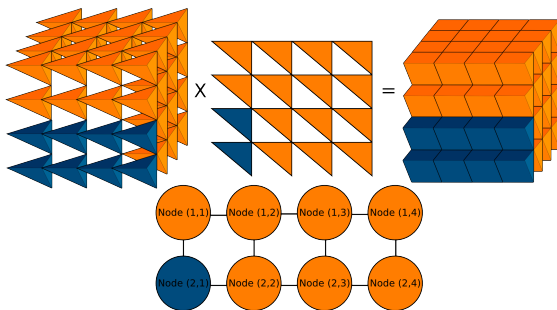
Parallel contraction in Cyclops

Cyclops uses nested parallel matrix multiplication variants

- 1D variants
 - perform a different *matrix-vector product* on each processor
 - perform a different *outer product* on each processor
- 2D variants
 - perform a different *inner product* on each processor
 - *scale a vector* on each processor then sum
- 3D variants
 - perform a different *scalar product* on each processor then sum
 - can be achieved by *nesting* 1D+1D+1D or 2D+1D or 1D+2D
- All variants are *blocked* in practice, naturally generalized to sparse matrix products

Tensor blocking/virtualization

Preserving symmetric-packed layout using cyclic distribution constrains possible tensor blockings



subdivision into more blocks than there are processors (virtualization)

Data mapping and redistribution

Transitions between contractions require redistribution and refolding

- 1D/2D/3D variants naturally map to 1D/2D/3D processor grids
- Initial tensor distribution is oblivious of contraction
 - by default each tensor distributed over all processors
 - user can specify [any processor grid mapping](#)
- Global redistribution done by one of three methods
 - reassign tensor blocks to processors (easy+fast)
 - reorder and reshuffle data to satisfy new blocking (fast)
 - treat tensors as sparse and sort globally by function of index
- Matricization/transposition is then done locally
 - dense tensor transpose done using [HPTT](#) (by Paul Springer)
 - sparse tensor converted to [CSR](#) sparse matrix format

Local summation and contraction

- For contractions, local summation and contraction is done via BLAS, including **batched GEMM**
- Threading is used via BLAS (done via OpenMP everywhere else)
- **GPU offloading** is available but not yet fully robust
- For sparse matrices, *MKL provides fast sparse matrix routines*
- To support **general (mixed-type, user-defined) elementwise functions**, manual implementations are available
- User can specify blocked implementation of their function to improve performance

- Performance models used to select best contraction algorithm
- Based on *linear cost model for each kernel*

$$T \approx \underbrace{\alpha S}_{\text{latency}} + \underbrace{\beta W}_{\text{comm. bandwidth}} + \underbrace{\nu Q}_{\text{mem. bandwidth}} + \underbrace{\gamma F}_{\text{flops}}$$

- Scaling of S , W , Q , F is a function of parameters of each kernel
- Coefficients for all kernels depend on compiler/architecture
- Linear regression with Tikhonov regularization used to select coefficients \mathbf{x}^*
- Model training done by benchmark suite that executes various end-functionality for growing problem sizes, collecting observations of parameters in rows of \mathbf{A} and execution timing in \mathbf{t}

$$\mathbf{x}^* = \underset{\mathbf{x}}{\operatorname{argmin}} (||\mathbf{A}\mathbf{x} - \mathbf{t}||_2 + \lambda ||\mathbf{x}||_2)$$

- Using Cython, we have provided a Python interface for Cyclops
- Follows `numpy.ndarray` conventions, plus sparsity and MPI execution

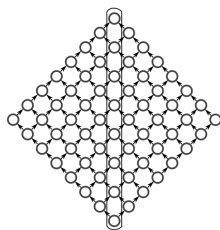
```
Z["abij"] += V["ijab"];           // C++  
Z.i("abij") << V.i("ijab")       // Python  
W["mnij"] += 0.5*W["mnef"]*T["efij"]; // C++  
W.i("mnij") << 0.5*W.i("mnef")*T.i("efij") // Python  
einsum("mnef,efij->mnij",W,T) // numpy-style Python
```

- Python interface is under active development, but is functional and available (**DEMO**)

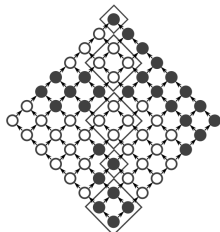
Tradeoffs in the diamond DAG

Computation vs synchronization tradeoff for the $n \times n$ diamond DAG,¹

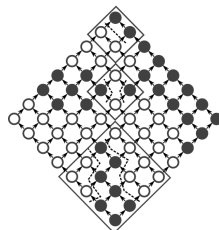
$$F \cdot S = \Omega(n^2)$$



Dependency chain P



Monochrome dependency intervals



Multicolored dependency intervals

In this DAG, vertices denote scalar computations in an algorithm

¹C.H. Papadimitriou, J.D. Ullman, SIAM JC, 1987

Scheduling tradeoffs of path-expander graphs

Definition ((ϵ, σ)-path-expander)

Graph $G = (V, E)$ is a (ϵ, σ) -**path-expander** if there exists a path $(u_1, \dots, u_n) \subset V$, such that the dependency interval $[u_i, u_{i+b}]_G$ for each i , b has size $\Theta(\sigma(b))$ and a minimum cut of size $\Omega(\epsilon(b))$.

Theorem (Path-expander communication lower bound)

*Any parallel schedule of an algorithm with a (ϵ, σ) -**path-expander** dependency graph about a path of length n and some $b \in [1, n]$ incurs computation (F), communication (W), and synchronization (S) costs:*

$$F = \Omega(\sigma(b) \cdot n/b), \quad W = \Omega(\epsilon(b) \cdot n/b), \quad S = \Omega(n/b).$$

Corollary

If $\sigma(b) = b^d$ and $\epsilon(b) = b^{d-1}$, the above theorem yields,

$$F \cdot S^{d-1} = \Omega(n^d), \quad W \cdot S^{d-2} = \Omega(n^{d-1}).$$

3D algorithms for matrix computations

For Cholesky factorization with p processors, cost is

$$F = \Theta(n^3/p), \quad W = \Theta(n^2/p^\delta), \quad S = \Theta(p^\delta)$$

for any $\delta = [1/2, 2/3]$.

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
SVD	singular values only \checkmark^5	singular vectors \times

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

¹B. Lipshitz, MS thesis 2013

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

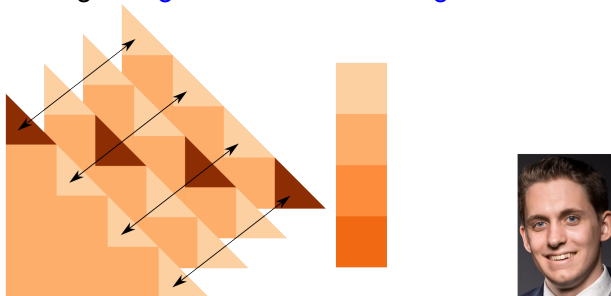
³A. Tiskin, FGCS 2007

⁴E.S., J. Demmel, EuroPar 2011

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

New algorithms can circumvent lower bounds

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



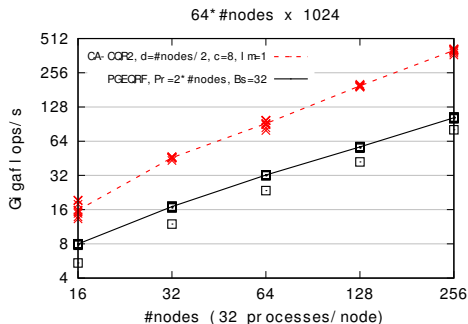
- **decreases synchronization cost** by $O(p^{2/3})$ on p processors with respect to known algorithms
- optimal communication for **any number of right-hand sides**
- MS thesis work by Tobias Wicky¹

¹T. Wicky, E.S., T. Hoefler, IPDPS 2017

Cholesky-QR2 for rectangular matrices

Cholesky-QR2¹ with 3D Cholesky gives a practical 3D QR algorithm

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



work by Edward Hutter (PhD student at UIUC)

¹Fukaya T, Nakatsukasa Y, Yanagisawa Y, Yamamoto Y. 2014

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

$$T = Q^T A Q$$

via a **two-sided orthogonal transformation** is most costly in diagonalization (eigenvalue computation, SVD similar)

- can be done by **successive column QR factorizations**

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

- 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

$$\underbrace{(I - UTU^T)^T}_{Q^T} A \underbrace{(I - UTU^T)}_Q = A - UV^T - VU^T$$

where columns of U are Householder vectors and V is

$$V^T = TU^T + \frac{1}{2}T^TU^T \underbrace{AU}_{\text{challenge}} TU^T$$

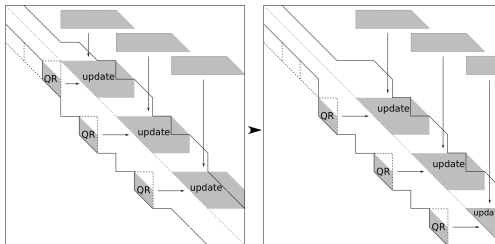
- if $b = 1$, U is a column-vector, and AU is dominated by **vertical communication cost** (moving A between memory and cache)
- **idea**: reduce to banded matrix ($b \gg 1$) first¹

¹ Auckenthaler, Bungartz, Huckle, Krämer, Lang, 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 band¹



- communication- and synchronization-efficient **1D SBR algorithm** known for small band-width²

¹ Lang 1993; Bischof, Lang, Sun 2000

² Ballard, Demmel, Knight 2012

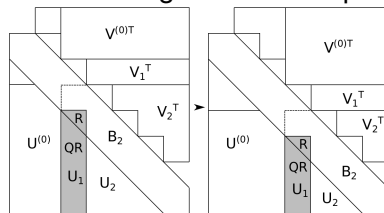
Communication-efficient eigenvalue computation

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

- implemented in ELPA, can outperform ScaLAPACK¹
- with $n = n/\sqrt{p}$, 1D SBR gives $W = O(n^2/\sqrt{p})$, $S = O(\sqrt{p} \log^2(p))^2$

New results³: **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)

¹ Auckenthaler, Bungartz, Huckle, Krämer, Lang, Willems 2011

² Ballard, Demmel, Knight 2012

³ S., Ballard, Demmel, Hoefler 2017

Symmetric eigensolver results summary

Algorithm	W	Q	S
ScaLAPACK	n^2/\sqrt{p}	n^3/p	$n \log(p)$
ELPA	n^2/\sqrt{p}	-	$n \log(p)$
two-stage + 1D-SBR	n^2/\sqrt{p}	$n^2 \log(n)/\sqrt{p}$	$\sqrt{p}(\log^2(p) + \log(n))$
many-stage	$n^2/p^{2/3}$	$n^2 \log(p)/p^{2/3}$	$p^{2/3} \log^2 p$

- 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)

Future directions and acknowledgements

Future/ongoing directions in Cyclops development

- General abstractions for tensor decompositions
- Concurrent scheduling of multiple contractions
- Fourier transforms along tensor modes
- Faster/specialized/optimized tensor slicing

Open problems in communication-avoiding matrix factorizations

- Algorithm for SVD with 3D cost without log factor overhead flops
- 3D implementation of SVD and QR with column pivoting

Acknowledgements

- Devin Matthews (UT Austin), Jeff Hammond (Intel Corp.), Grey Ballard, James Demmel (UC Berkeley), Tobias Wicky, Torsten Hoefler (ETH Zurich), Edward Hutter, Zecheng Zhang, Eric Song, Eduardo Yap, Linjian Ma (UIUC)
- Computational resources at NERSC, CSCS, ALCF, NCSA, TACC

Fast algorithms for symmetric tensor contractions

A tensor $\mathbf{T} \in \mathbb{R}^{n_1 \times \dots \times n_d}$ has

- order d (i.e. d modes / indices)
- dimensions n -by- \dots -by- n
- elements $T_{i_1 \dots i_d} = T_{\mathbf{i}}$ where $\mathbf{i} \in \{1, \dots, n\}^d$

We say a tensor is **symmetric** if for any $j, k \in \{1, \dots, n\}$

$$T_{i_1 \dots i_j \dots i_k \dots i_d} = T_{i_1 \dots i_k \dots i_j \dots i_d}$$

A tensor is **partially-symmetric** if such index interchanges are restricted to be within subsets of $\{1, \dots, n\}$, e.g.

$$T_{kl}^{ij} = T_{kl}^{ji} = T_{lk}^{ji} = T_{lk}^{ij}$$

For any $s, t, v \in \{0, 1, \dots\}$, a tensor contraction is

$$\forall \mathbf{i} \in \{1, \dots, n\}^s, \mathbf{j} \in \{1, \dots, n\}^t, \quad C_{ij} = \sum_{\mathbf{k} \in \{1, \dots, n\}^v} A_{i\mathbf{k}} B_{\mathbf{k}j}$$

Symmetric matrix times vector

Lets consider the simplest tensor contraction with symmetry

- let A be an n -by- n symmetric matrix ($A_{ij} = A_{ji}$)
- the symmetry is not preserved in matrix-vector multiplication

$$c = A \cdot b$$

$$c_i = \sum_{j=1}^n \underbrace{A_{ij} \cdot b_j}_{\text{nonsymmetric}}$$

- generally n^2 additions and n^2 multiplications are performed
- we can perform only $\binom{n+1}{2}$ multiplications using¹

$$c_i = \sum_{j=1, j \neq i}^n \underbrace{A_{ij} \cdot (b_i + b_j)}_{\text{symmetric}} + \underbrace{\left(A_{ii} - \sum_{j=1, j \neq i}^n A_{ij} \right) \cdot b_i}_{\text{low-order}}$$

¹E.S., J. Demmel, 2015

Symmetrized outer product

Consider a rank-2 outer product of vectors \mathbf{a} and \mathbf{b} of length n into symmetric matrix \mathbf{C}

$$\mathbf{C} = \mathbf{a} \cdot \mathbf{b}^\top + \mathbf{b} \cdot \mathbf{a}^\top$$
$$C_{ij} = \underbrace{\mathbf{a}_i \cdot \mathbf{b}_j}_{\text{nonsymmetric}} + \underbrace{\mathbf{a}_j \cdot \mathbf{b}_i}_{\text{permutation}}$$

usually computed via the n^2 multiplications and n^2 additions
new algorithm requires $\binom{n+1}{2}$ multiplications

$$C_{ij} = \underbrace{(\mathbf{a}_i + \mathbf{a}_j) \cdot (\mathbf{b}_i + \mathbf{b}_j)}_{\text{symmetric}} - \underbrace{\mathbf{a}_i \cdot \mathbf{b}_i}_{w_i} - \underbrace{\mathbf{a}_j \cdot \mathbf{b}_j}_{w_j}$$

z_{ij} w_i w_j

Symmetrized matrix multiplication

For symmetric matrices A and B , compute

$$C_{ij} = \sum_{k=1}^n \left(\underbrace{A_{ik} \cdot B_{kj}}_{\text{nonsymmetric}} + \underbrace{A_{jk} \cdot B_{ki}}_{\text{permutation}} \right)$$

New algorithm requires $\binom{n+2}{3} + \binom{n}{2}$ multiplications rather than n^3 , based on

$$\begin{aligned} C_{ij} &= \sum_{k=1}^n \left(\underbrace{A_{ik} \cdot B_{kj}}_{\text{nonsymmetric}} + \underbrace{A_{jk} \cdot B_{ki}}_{\text{transpose}} \right) \\ &= \sum_{k=1}^n \underbrace{(A_{ij} + A_{ik} + A_{jk}) \cdot (B_{ij} + B_{kj} + B_{ki})}_{\hat{Z}_{ijk}} - \underbrace{\sum_{k=1}^n A_{ik} \cdot B_{ik}}_{w_i} - \underbrace{\sum_{k=1}^n A_{jk} \cdot B_{jk}}_{w_j} \\ &\quad - nA_{ij} \cdot B_{ij} - A_{ij} \cdot \left(\underbrace{\sum_{k=1}^n B_{ki}}_{B_i^{(1)}} + \underbrace{\sum_{k=1}^n B_{kj}}_{B_j^{(1)}} \right) - \left(\underbrace{\sum_{k=1}^n A_{ki}}_{A_i^{(1)}} + \underbrace{\sum_{k=1}^n A_{kj}}_{A_j^{(1)}} \right) \cdot B_{ij}. \end{aligned}$$

Symmetrized tensor contraction

Generally consider any **symmetric tensor contraction** for $s, t, v \in \{0, 1, \dots\}$

$$\forall \mathbf{i} \in \{1, \dots, n\}^s, \mathbf{j} \in \{1, \dots, n\}^t, C_{ij} = \sum_{\mathbf{k} \in \{1, \dots, n\}^v} A_{ik} B_{kj} + \text{permutations}$$

- best previous algorithms used roughly $\binom{n}{s} \binom{n}{t} \binom{n}{v}$ multiplications, new algorithm requires roughly $\binom{n}{s+t+v}$ multiplications
- communication lower bounds for fast symmetric contractions
 - surprising **negative** result – when $s + t + v \geq 4$ and $s \neq t \neq v$ **asymptotically more communication** necessary for new algorithm!
- algorithm can be **nested** in the case of **partially-symmetric** contractions, leads to a **reduction in cost** – manyfold cost improvements in some high-order quantum chemistry methods