

CS 598 EVS: Tensor Computations

Bilinear Algorithms

Edgar Solomonik

University of Illinois at Urbana-Champaign

Bilinear Problems

- ▶ A number of basic numerical problems can be thought of as bilinear functions associated with particular order 3 tensors
 - ▶ *matrix multiplication*
 - ▶ *discrete convolution*
 - ▶ *symmetric tensor contractions*
- ▶ These problems admit nontrivial fast *bilinear algorithms*, which correspond to low-rank CP decompositions of the tensors
 - ▶ *Strassen's $O(n^{\log_2(7)})$ algorithm for matrix multiplication as well as all other subcubic matrix multiplication*
 - ▶ *The discrete Fourier transform (DFT), Toom-Cook, and Winograd algorithms for convolution are also examples of bilinear algorithms*
- ▶ *We will review fast bilinear algorithms for all of these approaches, using 0-based indexing when discussing convolution*

Bilinear Problems

- ▶ A bilinear problem for any inputs $\mathbf{a} \in \mathbb{R}^n$ and $\mathbf{b} \in \mathbb{R}^k$ computes $\mathbf{c} \in \mathbb{R}^m$ as defined by a tensor $\mathcal{T} \in \mathbb{R}^{m \times n \times k}$

$$c_i = \sum_{j,k} t_{ijk} a_j b_k \quad \Leftrightarrow \quad \mathbf{c} = \mathbf{f}^{(\mathcal{T})}(\mathbf{a}, \mathbf{b})$$

- ▶ Variants of discrete convolutions (linear convolution, correlation, cyclic convolution) provide simple examples of \mathcal{T}

- ▶ *Linear convolution*

$$t_{ijk} = \begin{cases} 1 & : k + i - j = 0 \\ 0 & : \text{otherwise} \end{cases} \quad \Rightarrow \quad c_i = \sum_{j,k} t_{ijk} a_j b_k = \sum_{j=\max(0, i-n+1)}^{\min(i, n-1)} a_j b_{i-j}$$

- ▶ *Correlation obtained by transposing the first and last mode of the linear convolution tensor*
- ▶ *Cyclic convolution has $t_{ijk} = 1$ if and only if $k + i - j = 0 \pmod{n}$*

Bilinear Algorithms

A bilinear algorithm (V. Pan, 1984) $\Lambda = (\mathbf{F}^{(A)}, \mathbf{F}^{(B)}, \mathbf{F}^{(C)})$ computes

$$\mathbf{c} = \mathbf{F}^{(C)}[(\mathbf{F}^{(A)T} \mathbf{a}) \odot (\mathbf{F}^{(B)T} \mathbf{b})],$$

where \mathbf{a} and \mathbf{b} are inputs and \odot is the Hadamard (pointwise) product.

$$\begin{bmatrix} \mathbf{c} \end{bmatrix} = \begin{bmatrix} \begin{matrix} \times & & \times & & \times & & \times \\ \times & \times & & \times & \times & & \times \\ \times & & \times & & \times & \times & \times \\ \times & \times & & & & \times & \times \\ \times & & \times & \times & & \times & \times \\ \times & \times & & \times & \times & \times & \times \\ \times & & \times & & \times & \times & \times \end{matrix} \end{bmatrix} \left[\left(\begin{bmatrix} \begin{matrix} \times \times & \times & \times & \times & \times & \times & \times \\ \times & \times & \times \times & \times & \times & \times & \times \\ \times & \times & & \times & \times \times & & \\ \times & \times & \times & & \times & \times \times & \\ \times & \times & \times & \times & \times & \times & \\ \times & & \times & \times & \times & \times & \\ \times & \times & \times & \times \times & \times & & \end{matrix} \end{bmatrix}^T \begin{bmatrix} \mathbf{a} \end{bmatrix} \right) \circ \left(\begin{bmatrix} \begin{matrix} \times & \times & \times & \times & \times & \times \\ \times & & \times & \times & \times & \\ \times \times & \times \times & \times & \times & \times \\ & \times \times & \times & \times & \\ \times & \times & \times & \times & \times \\ \times \times & \times & \times & \times \times \times \end{matrix} \end{bmatrix}^T \begin{bmatrix} \mathbf{b} \end{bmatrix} \right) \right]$$

Bilinear Algorithms as Tensor Factorizations

- ▶ A bilinear algorithm corresponds to a CP tensor decomposition

$$\begin{aligned}c_i &= \sum_{r=1}^R f_{ir}^{(C)} \left(\sum_j f_{jr}^{(A)} a_j \right) \left(\sum_k f_{kr}^{(B)} b_k \right) \\ &= \sum_j \sum_k \left(\sum_{r=1}^R f_{ir}^{(C)} f_{jr}^{(A)} f_{kr}^{(B)} \right) a_j b_k \\ &= \sum_j \sum_k t_{ijk} a_j b_k \quad \text{where} \quad t_{ijk} = \sum_{r=1}^R f_{ir}^{(C)} f_{jr}^{(A)} f_{kr}^{(B)}\end{aligned}$$

- ▶ For multiplication of $n \times n$ matrices, we can define a *matrix multiplication tensor* and consider algorithms with various bilinear rank
 - ▶ \mathbf{T} is $n^2 \times n^2 \times n^2$
 - ▶ Classical algorithm has rank $R = n^3$
 - ▶ Strassen's algorithm has rank $R \approx n^{\log_2(7)}$

Strassen's Algorithm

$$\text{Strassen's algorithm } \begin{bmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{bmatrix} = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \cdot \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix}$$

$$M_1 = (A_{11} + A_{22}) \cdot (B_{11} + B_{22})$$

$$C_{11} = M_1 + M_4 - M_5 + M_7$$

$$M_2 = (A_{21} + A_{22}) \cdot B_{11}$$

$$C_{21} = M_2 + M_4$$

$$M_3 = A_{11} \cdot (B_{12} - B_{22})$$

$$C_{12} = M_3 + M_5$$

$$M_4 = A_{22} \cdot (B_{21} - B_{11})$$

$$C_{22} = M_1 - M_2 + M_3 + M_6$$

$$M_5 = (A_{11} + A_{12}) \cdot B_{22}$$

$$M_6 = (A_{21} - A_{11}) \cdot (B_{11} + B_{12})$$

$$M_7 = (A_{12} - A_{22}) \cdot (B_{21} + B_{22})$$

By performing the nested calls recursively, Strassen's algorithm achieves cost,

$$T(n) = 7T(n/2) + O(n^2) = O(7^{\log_2 n}) = O(n^{\log_2 7})$$

Fast Bilinear Algorithms for Convolution

- ▶ Linear convolution corresponds to polynomial multiplication
 - ▶ *Let a and b be coefficients of degree $n - 1$ polynomial p and degree $k - 1$ polynomial q then*

$$(p \cdot q)(x) = \sum_{i=0}^{n+k-1} c_i x^i \quad \text{where} \quad c_i = \sum_{j=\max(0, i-n+1)}^{\min(i, n-1)} a_j b_{i-j}$$

- ▶ *This view motivates algorithms based on polynomial interpolation*
- ▶ The **Toom-Cook** convolution algorithm computes the coefficients of $p \cdot q$ by computing $(p \cdot q)(x_i)$ for $i \in \{1, \dots, n + k - 1\}$ and interpolates
 - ▶ *Let V_r be a $(n + k - 1)$ -by- r Vandermonde matrix based on the nodes x , so that $V_n \mathbf{a} = [p(x_1), \dots, p(x_{n+k-1})]^T$, etc.*
 - ▶ *Then to evaluate p and q at x and interpolate, we compute*

$$\mathbf{c} = V_{n+k-1}^{-1} ((V_n \mathbf{a}) \odot (V_k \mathbf{b}))$$

which is a bilinear algorithm

Toom-Cook Convolution and the Fourier Transform

- ▶ Vandermonde matrices are ill-conditioned with real nodes, but can be perfectly conditioned with complex nodes
 - ▶ *The condition number of a Vandermonde matrix with real nodes is exponential in its dimension*
 - ▶ *Choosing the nodes x to be the complex roots of unity gives the **discrete Fourier transform (DFT) matrix** $\mathbf{D}^{(n)}$, $d_{jk}^{(n)} = \omega_n^{jk}$ where $\omega_n = e^{2i\pi/n}$*
 - ▶ *Modulo normalization DFT matrix is orthogonal and symmetric (not Hermitian)*
- ▶ The **fast Fourier transform (FFT)** can be used to perform products with the DFT matrix in $O(n \log n)$ time Taking $\tilde{\mathbf{D}}^{(n)}$ to be the $n_1 \times n_2$ (for $n = n_1 n_2$) leading minor of \mathbf{D}_n we can compute $\mathbf{y} = \mathbf{D}^{(n)} \mathbf{x}$ via the split-radix- n_1 FFT,

$$y_k = \sum_{i=0}^{n-1} x_i \omega_n^{ik} = \sum_{i=0}^{n/2-1} x_{2i} \omega_{n/2}^{ik} + \omega_n^k \sum_{i=0}^{n/2-1} x_{2i+1} \omega_{n/2}^{ik}$$
$$y_{(kn_1+t)} = \sum_{s=0}^{n_1-1} \omega_{n_1}^{st} \left[\omega_n^{sk} \sum_{i=0}^{n_2-1} x_{(in_1+s)} \omega_{n_2}^{ik} \right] \Leftrightarrow \mathbf{Y} = ([\tilde{\mathbf{D}}^{(n)} \odot (\mathbf{D}^{(n_2)} \mathbf{A})] \mathbf{D}^{(n_1)})^T$$

Cyclic Convolution via DFT

- ▶ For linear convolution $D^{(n+k-1)}$ is used, for cyclic convolution $D^{(n)}$ suffices
 - ▶ Expanding the bilinear algorithm, $\mathbf{y} = D^{(n)^{-1}} ((D^{(n)} \mathbf{f}) \odot (D^{(n)} \mathbf{g}))$, we obtain

$$y_k = \frac{1}{n} \sum_{i=0}^{n-1} \omega_{(n)}^{-ki} \left(\sum_{j=0}^{n-1} \omega_{(n)}^{ij} f_j \right) \left(\sum_{t=0}^{n-1} \omega_{(n)}^{it} g_t \right) = \frac{1}{n} \sum_{i=0}^{n-1} \sum_{j=0}^{n-1} \sum_{t=0}^{n-1} \omega_{(n)}^{(j+t-k)i} f_j g_t$$

- ▶ It suffices to observe that for any fixed $u = j + t - k \neq 0$ or $\neq n$, the outer summation yields a zero result, since the geometric sum simplifies to

$$\sum_{i=0}^{n-1} \omega_{(n)}^{ui} = (1 - (\omega_{(n)}^u)^n) / (1 - \omega_{(n)}^u) = 0$$

- ▶ The DFT also arises in the eigendecomposition of a circulant matrix
 - ▶ The cyclic convolution is defined by the matrix-vector product $\mathbf{y} = C_{\langle \mathbf{a} \rangle} \mathbf{b}$ where

$$C_{\langle \mathbf{a} \rangle} = \begin{bmatrix} a_0 & \cdots & a_1 \\ \vdots & \ddots & \vdots \\ a_{n-1} & \cdots & a_0 \end{bmatrix}$$

- ▶ The eigenvalue decomposition of this matrix is $C_{\langle \mathbf{a} \rangle} = D^{(n)^{-1}} \text{diag}(D^{(n)} \mathbf{a}) D^{(n)}$

Winograd's Algorithm for Convolution

- ▶ The DFT/FFT requires complex arithmetic, motivating alternatives such as the more general Winograd family of algorithms
 - ▶ *In Winograd's convolution algorithm, the remainder of the product $v = pq$ is computed using k distinct polynomial divisors, $m^{(i)}$, whose product is the polynomial M with $\deg(M) > \deg(v)$*
 - ▶ *The k polynomial divisors, $m^{(1)}, m^{(2)}, \dots, m^{(k)}$ must be coprime*
 - ▶ *From the k remainders, $u^{(i)} = pq \bmod m^{(i)}$ the remainder $v = pq \bmod M$ is recovered via the Chinese remainder theorem*
 - ▶ *The theorem leverages Bézout's identity, which states that there exist polynomials $n^{(i)}$ and $N^{(i)}$ such that, for $M^{(i)} = M/m^{(i)}$,*

$$M^{(i)}N^{(i)} + m^{(i)}n^{(i)} = 1$$

which allow us to construct v

$$v = \left(\sum_{i=1}^k u^{(i)} M^{(i)} N^{(i)} \right) \bmod M$$

- ▶ *Toom-Cook algorithms are special cases of Winograd's convolution algorithm, where the polynomial divisors are $m^{(i)}(x) = x - \chi_i$, where χ_i are nodes*

Algebraic Formulation of Winograd's Algorithm for Convolution

- ▶ Given an operator $\mathbf{X}_{\langle m, d \rangle} \in \mathbb{C}^{\deg(m) \times (d+1)}$ to compute coefficients of $\rho = p \pmod{m}$, we can efficiently compute

$$pq \pmod{m} = (p \pmod{m})(q \pmod{m}) \pmod{m},$$

$$\mathbf{X}_{\langle m, \deg(p) + \deg(q) - 1 \rangle}(\mathbf{p} * \mathbf{q}) = \mathbf{X}_{\langle m, 2\deg(m) - 1 \rangle}((\mathbf{X}_{\langle m, \deg(p) \rangle} \mathbf{p}) * (\mathbf{X}_{\langle m, \deg(q) \rangle} \mathbf{q}))$$

- ▶ Further, given a bilinear algorithm $(\mathbf{A}, \mathbf{B}, \mathbf{C})$ to compute linear convolution of two m -dimensional vectors, we can obtain a bilinear algorithm $(\mathbf{X}_{\langle m, \deg(p) \rangle}^T \mathbf{A}, \mathbf{X}_{\langle m, \deg(q) \rangle}^T \mathbf{B}, \mathbf{X}_{\langle m, 2\deg(m) - 1 \rangle} \mathbf{C})$ to compute $\rho = pq \pmod{m}$, since

$$\rho = \mathbf{X}_{\langle m, 2\deg(m) - 1 \rangle} \mathbf{C}((\mathbf{A}^T \mathbf{X}_{\langle m, \deg(p) \rangle} \mathbf{p}) \odot (\mathbf{B}^T \mathbf{X}_{\langle m, \deg(q) \rangle} \mathbf{q})).$$

Algebraic Formulation of Winograd's Algorithm for Convolution

- ▶ Winograd's convolution algorithm effectively merges smaller bilinear algorithms for linear convolution
 - ▶ Given $M = \prod_{i=1}^k m^{(i)}$ where $\deg(M) = n + r - 1$ and $m^{(1)}, \dots, m^{(k)}$ are coprime, as well as $(A^{(i)}, B^{(i)}, C^{(i)})$ for $i \in \{1, \dots, k\}$, where $(A^{(i)}, B^{(i)}, C^{(i)})$ is a bilinear algorithm for linear convolution of vectors of dimension $\deg(m^{(i)})$
 - ▶ Winograd's convolution algorithm yields a bilinear algorithm (A, B, C) for computing linear convolution with vectors of dimension r and n , where

$$\begin{aligned} A &= \left[\mathbf{X}_{\langle m^{(1)}, r-1 \rangle}^T \mathbf{A}^{(1)} \quad \dots \quad \mathbf{X}_{\langle m^{(k)}, r-1 \rangle}^T \mathbf{A}^{(k)} \right], \\ B &= \left[\mathbf{X}_{\langle m^{(1)}, n-1 \rangle}^T \mathbf{B}^{(1)} \quad \dots \quad \mathbf{X}_{\langle m^{(k)}, n-1 \rangle}^T \mathbf{B}^{(k)} \right], \text{ and} \\ C &= \left[\tilde{C}^{(1)} \quad \dots \quad \tilde{C}^{(k)} \right] \end{aligned}$$

where $\tilde{C}^{(i)} = \mathbf{X}_{\langle M, \deg(M) + \deg(m^{(i)}) - 2 \rangle} \mathbf{T}_{\langle e^{(i)}, \deg(m^{(i)}) \rangle} \mathbf{X}_{\langle m^{(i)}, 2\deg(m^{(i)}) - 1 \rangle} C^{(i)}$ and $e^{(i)}$ are coefficients of polynomial $e^{(i)} = M^{(i)} N^{(i)} \bmod M$.

Algebraic Formulation of Winograd's Algorithm for Convolution

- ▶ A missing piece of the above formulation is how to realize Bézout's identity to compute $N^{(i)}$ and $e^{(i)}$
 - ▶ $e^{(i)} = M^{(i)}N^{(i)} \bmod M$ so it suffices to compute $n^{(i)}$ and $N^{(i)}$ then apply previously mentioned linear transformations
 - ▶ The extended Euclidian algorithm can be used for this task, or one can solve a linear system
 - ▶ The coefficients of polynomials \hat{N} and \hat{n} satisfying $\hat{M}\hat{N} + \hat{m}\hat{n} = 1$ for coprime \hat{M} and \hat{m} are

$$\begin{bmatrix} \hat{N} \\ \hat{n} \end{bmatrix} = \begin{bmatrix} \mathbf{T}_{\langle \hat{M}, \deg(\hat{m})-1 \rangle} & \mathbf{T}_{\langle \hat{m}, \deg(\hat{M})-1 \rangle} \end{bmatrix}^{-1} \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

Nested Bilinear Algorithms for Convolution

- ▶ 2D convolution is equivalent to nested 1D convolution
 - ▶ Given $F \in \mathbb{R}^{r \times r}$ and $G \in \mathbb{R}^{n \times n}$, the 2D linear convolution $Y = F * G$ with $Y \in \mathbb{R}^{(n+r-1) \times (n+r-1)}$ gives

$$y_{ab} = \sum_{i=\max(0, a-n+1)}^{\min(a, r-1)} \sum_{j=\max(0, b-n+1)}^{\min(b, r-1)} f_{ij} g_{a-i, b-j}$$

- ▶ 2D bilinear problem is defined by tensor $\mathcal{T}^{(2D)} = \mathcal{T} \otimes \mathcal{T}$ where \otimes is the natural generalization of Kronecker product to tensors
- ▶ 1D convolution can be reduced to 2D convolution with some work
 - ▶ For linear convolution, with vectors of dimension $n = st$ can reduce to $s \times t$ 2D convolution to obtain rank $(2s - 1)(2t - 1)$ bilinear algorithm via overlap-add technique, which computes partial sums of the result of the 2D convolution
 - ▶ For cyclic convolution, Agarwal-Cooley algorithm uses the Chinese remainder theorem for integers to decouple dimension $n = st$ convolution to $s \times t$ 2D cyclic convolution via permutations
- ▶ For more details on the above derivations and a broader survey of convolution algorithms, see <https://arxiv.org/abs/1910.13367>

Symmetric Tensor Contractions

- ▶ Bilinear algorithms can also be used to accelerate tensor contractions for tensors with symmetry
 - ▶ Recall a symmetric tensor is defined by e.g., $t_{ijk} = t_{ikj} = t_{kij} = t_{jki} = t_{jik} = t_{kji}$
 - ▶ Tensors can also have skew-symmetry (also known as antisymmetry, permutations have $+/-$ signs), partial symmetry (only some modes are permutable), or group symmetry (blocks are zero if indices satisfy modular equation)
 - ▶ The simplest example of a symmetric tensor contraction is

$$y = Ax \quad \text{where } A = A^T$$

it is not obvious how to leverage symmetry to reduce cost of this contraction

- ▶ Bilinear algorithms for symmetric tensor contractions exist with lower rank than their nonsymmetric counterparts
 - ▶ Symmetric matrix-vector product can be done with $n(n+1)/2$ multiplications
 - ▶ Cost of contractions of partially symmetric tensors reduced via this technique

Symmetric Matrix Vector Product

- ▶ Consider computing $c = \mathbf{A}b$ with $\mathbf{A} = \mathbf{A}^T$
 - ▶ Typically requires n^2 multiplications since $a_{ij}b_j \neq a_{ji}b_i$ and $n^2 - n$ additions
 - ▶ Instead can compute

$$v_i = \sum_{j=1}^{i-1} u_{ij} + \sum_{j=i+1}^n u_{ji} \quad \text{where} \quad u_{ij} = a_{ij}(b_i + b_j)$$

using $n(n-1)/2$ multiplications (since we only need u_{ij} for $i > j$) and about $3n^2/2$ additions, then

$$c_i = (2a_{ii} - \sum_{j=1}^n a_{ij})b_i + v_i$$

using n more multiplications and n^2 additions

- ▶ Beneficial when multiplying elements of \mathbf{A} and b costs more than addition
- ▶ This technique yields a bilinear algorithm with rank $n(n+1)/2$

Partially-Symmetric Tensor Times Matrix (TTM)

- ▶ Can use symmetric mat-vec algorithm to accelerate TTM with partially symmetric tensor from $2n^4$ operations to $(3/2)n^4 + O(n^3)$
 - ▶ Given $\mathcal{A} \in \mathbb{R}^{n \times n \times n}$ with symmetry $a_{ijk} = a_{jik}$ and $\mathbf{B} \in \mathbb{R}^{n \times n}$, we compute

$$c_{ikl} = \sum_j a_{ijk} b_{jl}$$

- ▶ We can think of this as a set of symmetric matrix-vector products

$$\mathbf{c}^{(k,l)} = \mathbf{A}^{(k)} \mathbf{b}^{(l)}$$

and apply the fast bilinear algorithm

$$v_{ikl} = \sum_{j=1}^{i-1} u_{ijkl} + \sum_{j=i+1}^n u_{ijkl} \quad \text{where} \quad u_{ijkl} = a_{ijk}(b_{il} + b_{jl})$$

$$c_{ikl} = (2a_{iik} - \sum_{j=1}^n a_{ijk})b_{il} + v_{ikl}$$

using about $n^4/2$ multiplications and $n^4 + O(n^3)$ additions (need only n^3 distinct sums of elements of \mathbf{B}) to compute \mathcal{V} , then $O(n^3)$ operations to get \mathcal{C} from \mathcal{V}

Computing Symmetric Matrices

- ▶ Output symmetry can also be used to reduced cost, for example when computing a symmetrized outer product $C = ab^T + ba^T$
 - ▶ $C = C^T$ so suffices to compute c_{ij} for $i \geq j$, $c_{ij} = a_i b_j + a_j b_i$
 - ▶ To reduce number of products by a factor of 2, can instead compute

$$c_{ij} = (a_i + a_j)(b_i + b_j) - v_i - v_j \quad \text{where} \quad v_i = a_i b_i$$

- ▶ To symmetrize product of two symmetric matrices, can compute anticommutator, $C = AB + BA$
 - ▶ Each matrix can be represented with $n(n+1)/2$ elements, but products all n^3 products $a_{ik} b_{kj}$ are distinct (so typically cost is $2n^3$)
 - ▶ Cost can be reduced to $n^3/6 + O(n^2)$ products by amortizing terms in

$$c_{ij} = \sum_k (a_{ij} + a_{ik} + a_{jk})(b_{ij} + b_{ik} + b_{jk}) - na_{ij}b_{ij} \\ - \left(\sum_k a_{ik} + a_{jk} \right) b_{ij} - a_{ij} \left(\sum_k b_{ik} + b_{jk} \right) - \sum_k a_{ik} b_{ik} - \sum_k a_{jk} b_{jk}$$

General Symmetric Tensor Contractions

- ▶ We can now consider the cost of a symmetrized contraction over v indices of symmetric tensors \mathcal{A} (of order $s + v$) and \mathcal{B} (of order $v + t$)

$$c_{i'_1 \dots i'_s, j'_1 \dots j'_t} = \sum_{\{i_1 \dots i_s, j_1 \dots j_t\} \in \Pi(i'_1 \dots i'_s, j'_1 \dots j'_t)} \sum_{k_1 \dots k_v} a_{i_1 \dots i_s, k_1 \dots k_v} b_{k_1 \dots k_v, j_1 \dots j_t}$$

where Π gives all distinct partitions of the $s + t$ indices into two subsets of size s and t , e.g.,

$$\Pi(i_1, j_1 j_2) = \{\{i_1, j_1 j_2\}, \{j_1, i_1 j_2\}, \{j_2, i_1 j_1\}\}$$

- ▶ Such tensor contractions can be done using $n^{s+t+v} / (s + t + v)! + O(n^{s+t+v-1})$ products
 - ▶ *General algorithm looks similar to anticommutator matrix product*
 - ▶ *After multiplying subsets of operands, unneeded terms are all computable with $O(n^{s+t+v-1})$ products*
 - ▶ *These approaches correspond to bilinear algorithms of this rank*

Hankel Matrix Vector Product

- ▶ A Hankel matrix is a reflection of a Toeplitz matrix (which can similarly be used to compute convolution),

$$\mathbf{H} = \begin{bmatrix} h_1 & h_2 & h_3 & \dots \\ h_2 & h_3 & & \\ h_3 & & \ddots & \\ \vdots & & & \end{bmatrix}$$

- ▶ *Hankel matrices are symmetric across blocks and Hankel within blocks*

$$\mathbf{H} = \begin{bmatrix} \mathbf{H}_1 & \mathbf{H}_2 \\ \mathbf{H}_2 & \mathbf{H}_1 \end{bmatrix}$$

but not vice versa, i.e., some elements in \mathbf{H}_1 must be the same as in \mathbf{H}_2

- ▶ Can compute Hankel mat-vec by applying symmetric mat-vec recursively
 - ▶ *For $n = 2$ fast bilinear algorithm requires $n(n + 1)/2 = 3$ products*
 - ▶ *Additions of Hankel matrices and vector entries can be done in $O(n)$ time, so*

$$T(n) = 3T(n/2) + O(n) = O(n^{\log_2(3)})$$

- ▶ *This complexity is as good as the recursive application of Toom-Cook, but $O(n \log n)$ can be achieved via FFT or other techniques*

Group Symmetry

- ▶ Tensors arising in physical simulations often have group structure that reflects conservation laws
 - ▶ *Abelian group symmetries can be mapped to cyclic group, which can be used to define a block-sparse form of the tensors (here represented using extra modes)*
 - ▶ *A particularly common/important contraction with cyclic group symmetry is*

$$w_{aA,bB,iI,jJ} = \sum_{k,K,l,L} u_{aA,bB,kK,lL} v_{kK,lL,iI,jJ}.$$

where for some group size G , we have symmetries, e.g.,

$$w_{aA,bB,iI,jJ} \neq 0 \text{ if } A + B - I - J \equiv 0 \pmod{G},$$

$$u_{aA,bB,kK,lL} \neq 0 \text{ if } A + B + K + L \equiv 0 \pmod{G},$$

$$v_{kK,lL,iI,jJ} \neq 0 \text{ if } K + L - I - J \equiv 0 \pmod{G}.$$

- ▶ *We can write each of these tensors using a **reduced form** and an **irrep map**,*

$$w_{aA,bB,iI,jJ} = r_{aA,bB,iI,jJ}^{(W)} m_{ABIJ}$$

where $m_{ABIJ} = 1$ if $A + B - I - J \equiv 0 \pmod{G}$ and $m_{ABIJ} = 0$ otherwise

Fast Algorithms for Contraction with Group Symmetry

- ▶ The irrep map tensor describes the cyclic group and has a lot of structure

- ▶ *From definition, $m_{ABIJ} = 1$ if $A + B - I - J \equiv 0 \pmod{G}$ and $m_{ABIJ} = 0$ otherwise, we see there are $O(G^3)$ nonzeros*
- ▶ *No matter the order of \mathcal{M} or the ordering of indices, it will have a tensor train decomposition of rank G , in particular,*

$$m_{ABIJ} = m_{ABQ}^{(1)} m_{QIJ}^{(2)} \quad \text{where} \quad m_{ABQ}^{(1)} = 1 \quad \text{iff} \quad A + B \equiv Q \pmod{G}, \quad \text{etc.,}$$

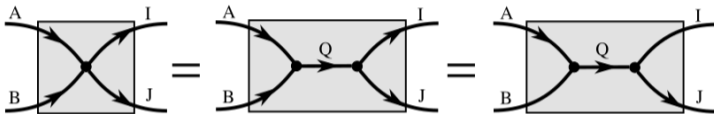
- ▶ *We can then also define a reduced form relative to an auxiliary index in such a decomposition, e.g.,*

$$w_{aA,bB,iI,jJ} = \sum_Q r_{aA,bB,iQ,j}^{(W)} m_{ABQ}^{(1)} m_{QIJ}^{(2)}$$

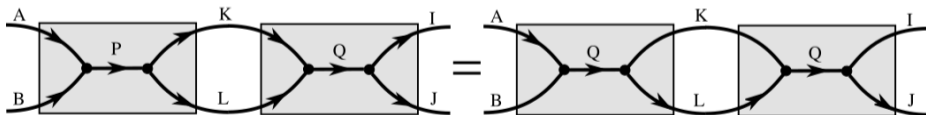
- ▶ *By defining consistent reduced forms using auxiliary indices, we can efficiently contraction group symmetric tensors via set of dense symmetric contractions*

Fast Algorithms for Contraction with Group Symmetry

- ▶ We can represent the special reduced form via a tensor diagram

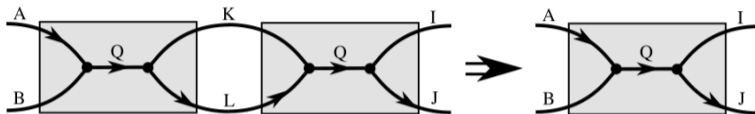


- ▶ Given a consistent choice of intermediate index in the two operand tensors, they can be unified



Fast Algorithms for Contraction with Group Symmetry

- ▶ After the indices are matched, we can contract efficiently



$$r_{aA,b,i,jJ,Q}^{(W)} = \sum_{k,l,L} r_{aA,b,k,lL,Q}^{(U)} r_{k,lL,i,jJ,Q}^{(V)}$$

- ▶ Naive contraction of original tensors had cost $O(n^6)$
 - ▶ *New algorithm has cost $O(n^6/G^2)$*
 - ▶ *Factor of G^2 improvement attainable for also for other (higher order) contractions*

Bilinear Algorithm for Contraction with Group Symmetry

- ▶ A group symmetric contraction is a bilinear algorithm
 - ▶ *Can view the contraction in a nested fashion as*

$$\mathcal{K}_{ABIJ}^{(W)} = \sum_{KL} \mathcal{K}_{ABKL}^{(U)} \cdot \mathcal{K}_{KLIJ}^{(V)}$$

where $\mathcal{F} = \mathcal{K}_{ABKL}^{(U)} \in \mathbb{R}^{n \times n \times n \times n}$ and $\mathcal{G} = \mathcal{K}_{KLIJ}^{(V)} \in \mathbb{R}^{n \times n \times n \times n}$, so that $\mathcal{H} = \mathcal{F} \cdot \mathcal{G}$ gives

$$h_{abij} = \sum_{kl} f_{abkl} g_{klij} = \sum_{kl} u_{aA,bB,kK,lL} v_{kK,lL,iI,jJ}$$

- ▶ *Therefore, it suffices to find a bilinear algorithm for the first contraction of $G \times G \times G \times G$ tensors*

$$k_{ABIJ}^{(W)} = \sum_{KL} k_{ABKL}^{(U)} \cdot k_{KLIJ}^{(V)}$$

with the same sparsity as irrep maps $\mathcal{M}^{(U)}$, $\mathcal{M}^{(V)}$, and $\mathcal{M}^{(W)}$

- ▶ *An algorithm for the full contraction can then be constructed by nesting with a standard blockwise contraction*

Bilinear Algorithm for Contraction with Group Symmetry

- ▶ The group symmetric contraction algorithm leverages a CP decomposition
 - ▶ *The tensor defining the bilinear problem group symmetric contraction can be written as*

$$T_{ABIJ, \hat{A}\hat{B}KL, \hat{K}\hat{L}\hat{I}\hat{J}} = \delta_{A+B, I+J} \delta_{A+B, -K-L} \prod_{X \in \{A, B, I, J, K, L\}} \delta(X, \hat{X})$$

- ▶ *Using the CP decomposition obtained via the identity $\delta_{A+B, I+J} \delta_{A+B, -K-L} = \sum_{Q=1}^G \delta_{AB, Q} \delta_{IJ, Q} \delta_{-K-L, Q}$, we can define a bilinear algorithm of rank G^4 that acts as follows*

$$k_{ABIJ}^{(W)} = \sum_{\hat{A}\hat{J}\hat{L}Q} \delta_{A, \hat{A}} \delta_{J, \hat{J}} \delta_{AB, Q} \delta_{IJ, Q} \left(\sum_{\hat{B}\hat{K}} \delta_{-K-L, Q} k_{\hat{A}\hat{B}KL}^{(U)} \right) \left(\sum_{\hat{I}\hat{K}} \delta_{\hat{I}\hat{J}, Q} k_{\hat{I}\hat{J}\hat{K}L}^{(V)} \right)$$