



BERKELEY LAB

Accelerating Multilinear Maps and Structured Sparse Tensor Kernels

EECS Department Dissertation Talk

July 21, 2025

Vivek Bharadwaj

Committee in Charge

Professor James Demmel, co-chair
Adj. Associate Professor Aydin Buluç, co-chair
Professor Katherine Yelick
Assistant Professor Michael Lindsey



Multilinear Algebra in the Exascale Era



- Linear algebra drives progress in scientific computing and machine learning.
- LLNL El Capitan: 1.7×10^{18} FP64 FLOPs on LINPACK, more in lower precision.
- Growing number of applications rely on *multilinear algebra*: much less studied / optimized. Let's survey three such applications.



Symmetry-Preserving Chemistry Models

- Scientists today can use **neural networks** to quickly approximate atomic forces, predict properties of molecular systems.
- State-of-art models even bake physical laws, like **symmetry**, into network structure.
- Accurate, high throughput simulation could revolutionize material discovery and drug design.

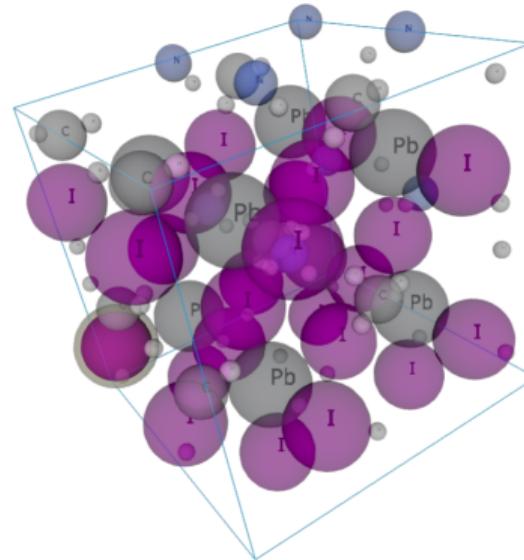
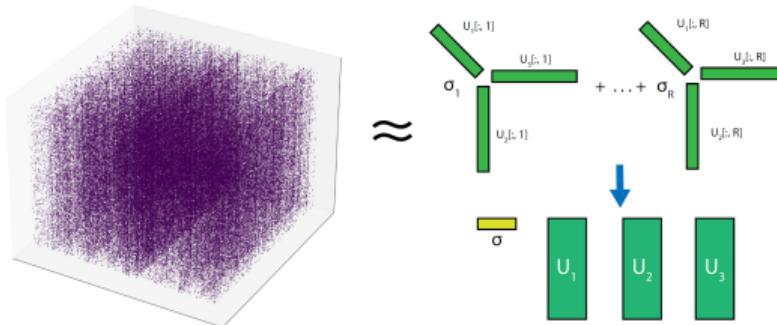


Figure: Perovskite crystal visualized with x3dase.

Multidimensional Network Analysis



- Consider a large, sparse array containing, e.g., e-commerce review data, Reddit word usage, etc.
- We have mathematical tools to cluster / analyze quantities of interest using **only the location and values** of nonzero elements by generalizing matrix PCA.



Inverse Problems and Tomography

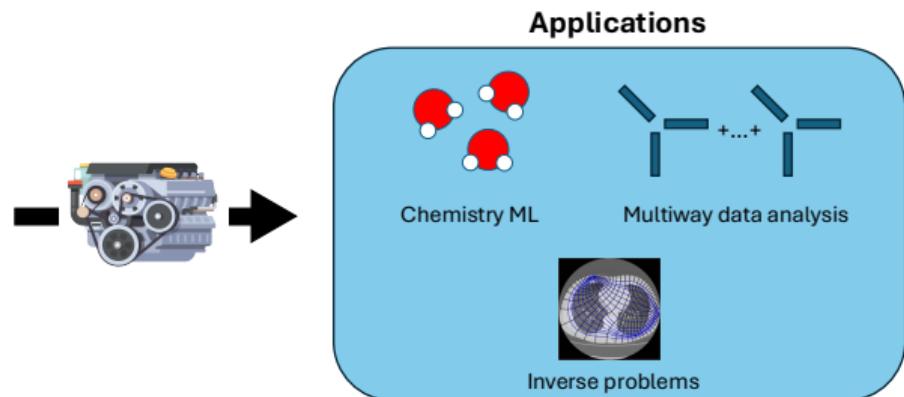
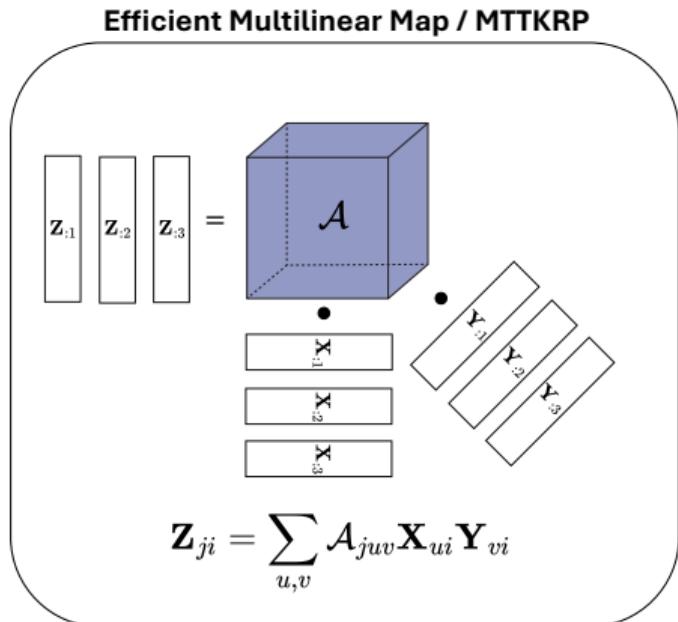
We can image interior of biological tissue by applying weak currents, measuring voltages on a boundary, solving linear least-squares problems [Che+20].



Figure: EIT demo on orange (CC2.0 Wikimedia, John Cummings).



Multilinear Maps: The Common Pattern





Warm-up: Linear Maps

- Consider $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ satisfying

$$f(\alpha x) = \alpha f(x) \quad f(x + y) = f(x) + f(y)$$

What linear algebraic primitive lets us evaluate such maps efficiently?

- Answer: Matrix-vector multiply for one input, matrix-matrix multiply enables high performance for multiple inputs: $f(X) = AX$, $A \in \mathbb{R}^{m \times n}$.

$$\begin{matrix} z \\ \end{matrix} = \begin{matrix} A \\ \end{matrix} \begin{matrix} x \\ \end{matrix}$$



Multilinear Maps and the MTTKRP

$$\begin{matrix} \textbf{z}_1 & \textbf{z}_2 & \textbf{z}_3 \end{matrix} = \textbf{A} \cdot \begin{matrix} \textbf{x}_1 & \textbf{x}_2 & \textbf{x}_3 \\ \otimes & \otimes & \otimes \\ \textbf{y}_1 & \textbf{y}_2 & \textbf{y}_3 \end{matrix}$$

$\textbf{Z}_{:i} = \textbf{A} \cdot (\textbf{X}_{:i} \otimes \textbf{Y}_{:i})$

$\textbf{Z}_{ji} = \sum_{u,v} \mathcal{A}_{juv} \textbf{X}_{ui} \textbf{Y}_{vi}$

$$\begin{matrix} \textbf{z}_1 & \textbf{z}_2 & \textbf{z}_3 \end{matrix} = \textbf{A} \cdot \begin{matrix} \textbf{x}_1 \\ \textbf{x}_2 \\ \textbf{x}_3 \end{matrix} \odot \begin{matrix} \textbf{y}_1 \\ \textbf{y}_2 \\ \textbf{y}_3 \end{matrix}$$

- Consider a multilinear function $f(\mathbf{x}, \mathbf{y})$. What fundamental kernels allow efficient evaluation?
- Answer: matrix-vector multiply by a **Kronecker product**:

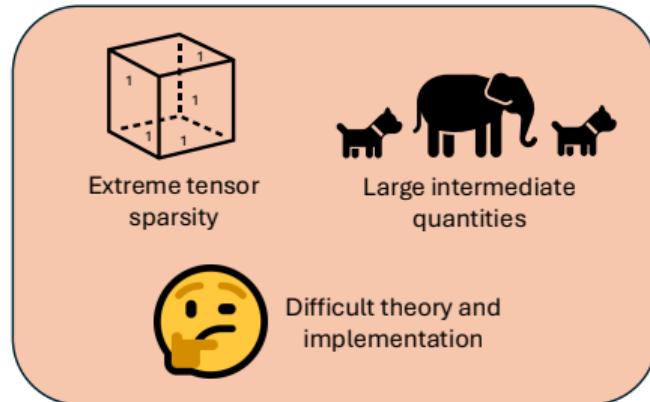
$$f(\mathbf{x}, \mathbf{y}) = \mathbf{A}(\mathbf{x} \otimes \mathbf{y}), \mathbf{A} \in \mathbb{R}^{k \times mn},$$

matrix-matrix multiplication by a **Khatri-Rao product** for a batch:

$$f(\mathbf{X}, \mathbf{Y}) = \mathbf{A}(\mathbf{X} \odot \mathbf{Y}).$$



Challenges and Opportunities



- Sparse tensors (e.g. Reddit-2015 [Smi+17]) exhibit nonzero fraction as low as 4×10^{-10} ; curse of dimensionality.
- Kronecker product height is exponential in input dimension; costly to materialize.
- Existing algorithms / software target matrices, not tensors.

We accelerate the multilinear map in two important applications, offering OOM speedups via kernel engineering and asymptotically faster approximation algorithms.

Research Papers Covered Today



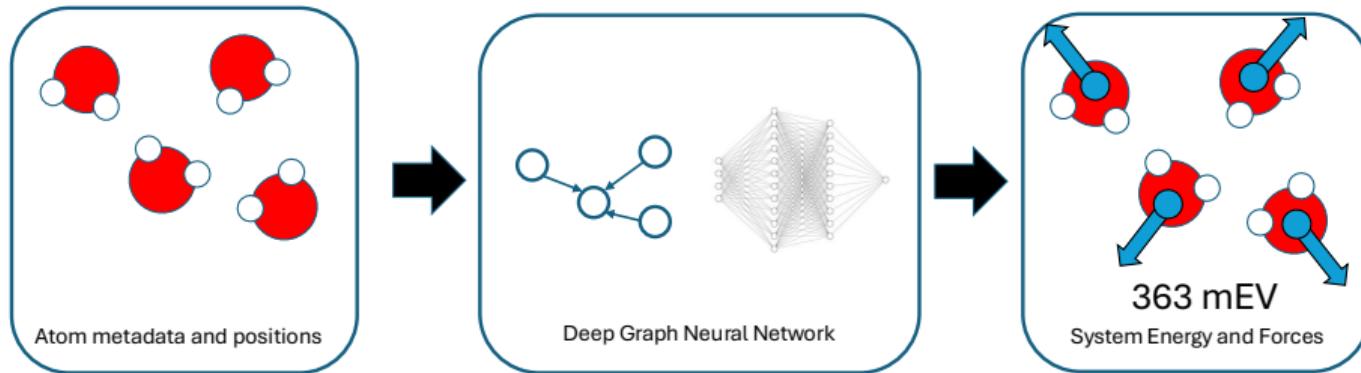
In this talk, we explore multiple methods to accelerate the **batched multilinear map**:

1. **Optimizing Chemistry Foundation Models:** An Efficient Sparse Kernel Generator for $O(3)$ -Equivariant Deep Networks. ACDA [Bha+25].
2. **Accelerated Randomized Tensor Decomposition:** Fast Exact Leverage Score Sampling from Khatri-Rao Products. NeurIPS [Bha+23].
3. Distributed-Memory Randomized Algorithms for Sparse Tensor CP Decomposition. SPAA [Bha+24a].
4. Efficient Leverage Score Sampling for Tensor Trains. NeurIPS [Bha+24b].
5. Distributed-Memory Sparse Kernels for Machine Learning. IPDPS [BBD22].



Part 1: Clebsch-Gordon Tensor Products

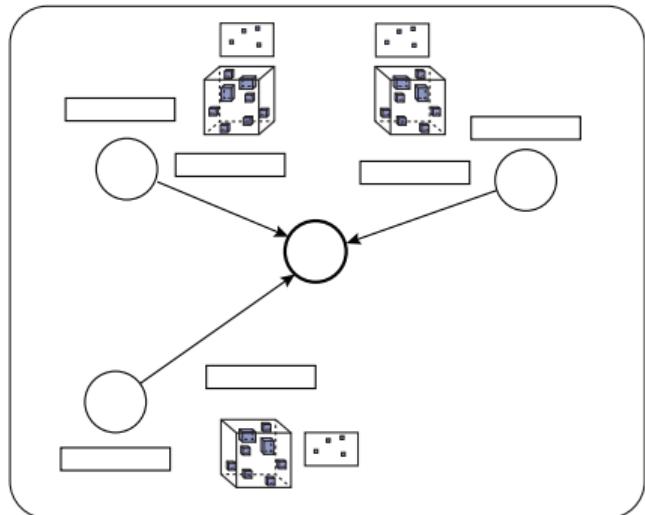
Geometric Deep Learning in Atomic Simulation



- Geometric deep learning achieves SOTA performance for fast interatomic potential calculation. Key primitive is the **message-passing graph neural network**.
- Input: atom metadata and positions. Output: system energy, atomic forces.



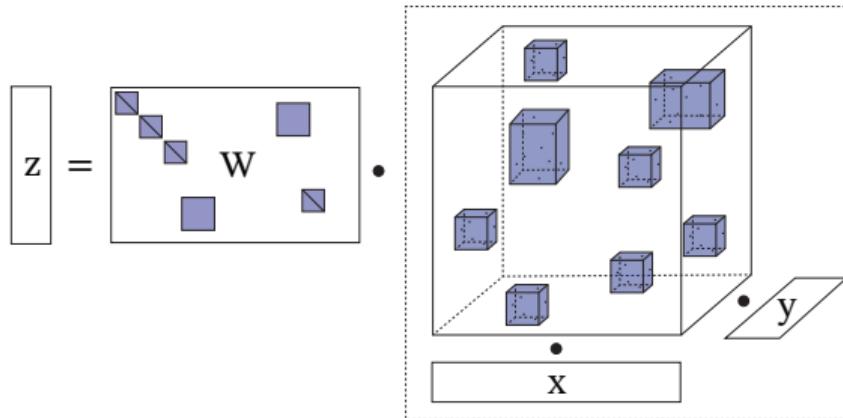
Respecting Physical Symmetries



- Message-passing GNNs generate messages for each node and edge.
- $O(3)$ -Equivariance: if the input coordinate system rotates, the output energy stays **the same** and the predicted forces **rotate compatibly**.
- Need to combine node, edge features in a highly structured, prescriptive manner.



The Clebsch-Gordon Tensor Product



$$z = W \cdot P \cdot (x \otimes y) = W \sum_{i=1, j=1}^{m, n} x[i] y[j] \mathcal{P}[ij :]$$

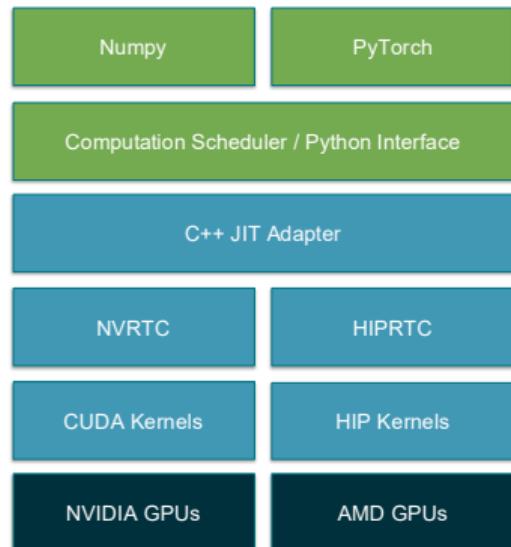
Kernel typically executes on a large batch of (x, y, W) inputs ($B \approx 10^5 - 10^7$).



Our Contributions [Bha+25]

We introduce *OpenEquivariance*, a fast kernel generator for CG tensor products with up to **10x speedup** over the popular e3nn package, on par with NVIDIA cuE (joint work with Austin Glover).

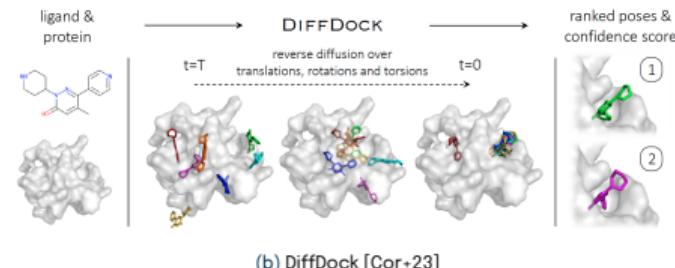
- Exploits sparse structure in the multilinear map tensor via JIT, register caching, loop unrolling.
- Provides FlashAttention-style [Dao+22] backward pass, novel identities for higher gradients.
- Fuses CG tensor product with graph convolution, saving orders of magnitude of memory.



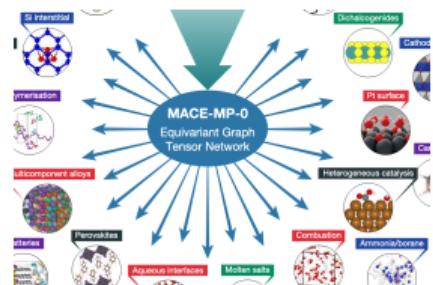
Applications of OpenEquivariance



(a) Nequip [Bat+22]



(b) DiffDock [Cor+23]

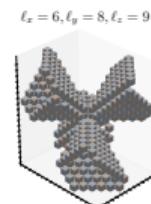
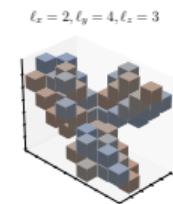
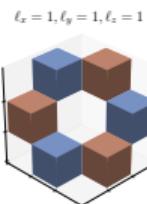


(c) MACE-MPO [Bat+24]



Subkernels of the CG Tensor Product

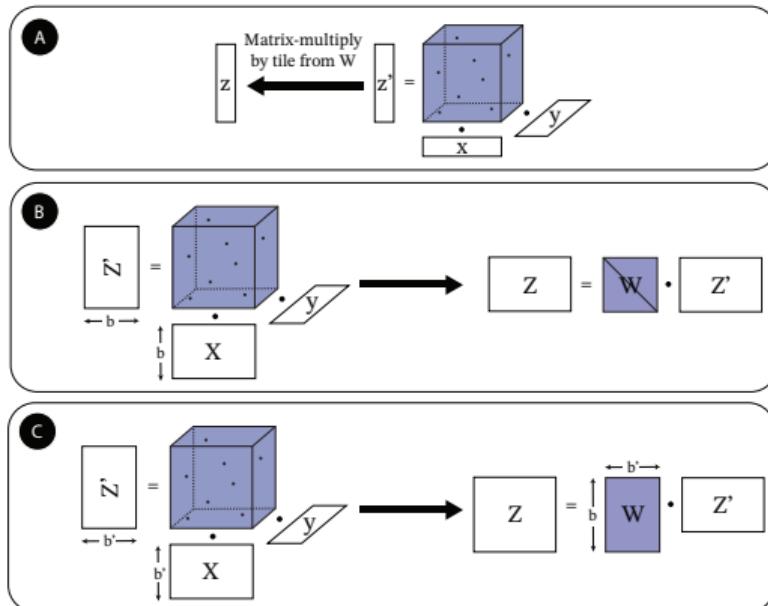
- Many nonzero blocks repeated in the sparse tensor (known at model compile-time).
- Operation splits into many smaller operations (subkernels)



nonzero % = 0.22

nonzero % = 0.16

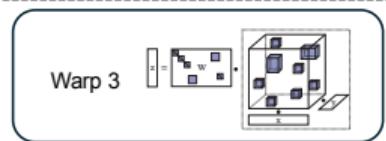
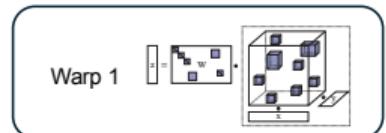
nonzero % = 0.08





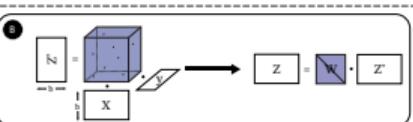
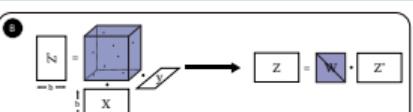
A Roadmap to Efficient CG Kernels

Assign each batch element to distinct GPU warp

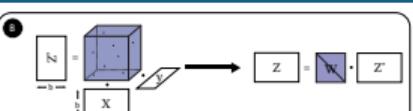


Schedule subkernels to manage SMEM usage, avoid memory traffic

Load subsegments of x , y , W into SMEM, reshape



Flush z to global memory, load next segment of x



Register-cache operands, use JIT to emit optimized instruction stream

Source

```
z[1] += 0x1.6a09e60000000p-2f * x[0] * y[0];
z[0] += -0x1.6a09e60000000p-2f * x[1] * y[0];
z[2] += 0x1.d363d00000000p-2f * x[1] * y[0];
z[1] += -0x1.d363d00000000p-2f * x[2] * y[0];
z[4] += -0x1.6a09e60000000p-1f * x[3] * y[0];
z[3] += 0x1.6a09e60000000p-1f * x[4] * y[0];
z[5] += -0x1.d363d00000000p-2f * x[4] * y[0];
z[4] += 0x1.d363d00000000p-2f * x[5] * y[0];
z[6] += -0x1.6a09e60000000p-2f * x[5] * y[0];
z[5] += 0x1.279a740000000p-1f * x[1] * y[1];
z[4] += 0x1.279a740000000p-2f * x[2] * y[1];
z[2] += -0x1.279a740000000p-2f * x[4] * y[1];
z[1] += -0x1.279a740000000p-1f * x[5] * y[1];
```

PTX

```
fma.rn.f32 %f4183, %f4126, %f4089, %f4182;
fma.rn.f32 %f4184, %f4152, %f4089, %f4183;
fma.rn.f32 %f4185, %f4112, %f4089, %f4184;
mul.f32 %f4186, %f4094, @f3f080677;
mul.f32 %f4187, %f4186, %f4096;
fma.rn.f32 %f4188, %f4187, %f4089, %f4146;
fma.rn.f32 %f4189, %f4100, %f4089, %f4188;
mul.f32 %f4190, %f4133, @f3EA79762;
mul.f32 %f4191, %f4190, %f4096;
fma.rn.f32 %f4192, %f4191, %f4089, %f4171;
fma.rn.f32 %f4193, %f4137, %f4089, %f4192;
fma.rn.f32 %f4194, %f4135, %f4089, %f4158;
fma.rn.f32 %f4195, %f4137, %f4089, %f4194;
```



Warp-Level Forward Algorithm

Algorithm Subkernel C Warp-Level Algorithm

Require: $\mathbf{X} \in \mathbb{R}^{b' \times (2\ell_x + 1)}$, $\mathbf{y} \in \mathbb{R}^{(2\ell_y + 1)}$, $\mathbf{W} \in \mathbb{R}^{b \times b'}$

Require: Sparse tensor $\mathcal{P}^{(\ell_x, \ell_y, \ell_z)}$ for subkernel

for $t = 1 \dots b'$ **do** ▷ Parallel over threads

 Load $\mathbf{x}_{\text{reg}} = \mathbf{X} [t, :]$, $\mathbf{y}_{\text{reg}} = \mathbf{y}$

 Initialize $\mathbf{z}_{\text{reg}} \in \mathbb{R}^{2\ell_z + 1}$ to 0.

for $(i, j, k, v) \in \text{nz}(\mathcal{P}^{(\ell_x, \ell_y, \ell_z)})$ **do** ▷ Unroll via JIT

$\mathbf{z}_{\text{reg}} [k] += v \cdot \mathbf{x}_{\text{reg}} [i] \cdot \mathbf{y}_{\text{reg}} [j]$

Store $\mathbf{Z}' [t, :] = \mathbf{z}_{\text{reg}}$, compute $\mathbf{Z} += \mathbf{W} \cdot \mathbf{Z}'$.

- Gradient calculation (backward pass) is similar, but requires 3 update equations instead of 1.

Optimizing Model Force Predictions



- The derivative of predicted energy $E(\mathbf{R}, \mathbf{W})$ w.r.t. atomic positions \mathbf{R} yields atomic forces. Want to solve

$$\min_{\mathbf{W}} \mathcal{L}(\mathbf{R}, \mathbf{W}) = \min_{\mathbf{W}} \|\mathbf{F}_{\text{pr}}(\mathbf{R}, \mathbf{W}) - \mathbf{F}_{\text{gt}}(\mathbf{R})\|_F^2$$

- To minimize difference between predicted forces and ground truth, need *second partial derivative of energy*:

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial \mathbf{W}} &= 2 \cdot \text{vec}(\mathbf{F}_{\text{pr}}(\mathbf{R}, \mathbf{W}) - \mathbf{F}_{\text{gt}}(\mathbf{R}))^\top \frac{\partial \mathbf{F}_{\text{pr}}}{\partial \mathbf{W}} \\ &= -2 \cdot \text{vec}(\mathbf{F}_{\text{pr}}(\mathbf{R}, \mathbf{W}) - \mathbf{F}_{\text{gt}}(\mathbf{R}))^\top \frac{\partial^2 E}{\partial \mathbf{R} \partial \mathbf{W}}\end{aligned}$$

Novel Identities for Second Partial Derivatives



- PyTorch gladly takes the second partial derivative... if you supply a “double-backward” pass for the CG tensor product.
- **Key:** No additional kernel engineering required! Double-backward can be expressed as a linear combination of forward / backward calls:

$$\text{op1} = \text{backward}(\partial \mathcal{L} / \partial \mathbf{a}, \partial \mathcal{L} / \partial \mathbf{b}, \mathbf{W}, \mathbf{g}_z)$$

$$\text{op2} = \text{backward}(\mathbf{x}, \mathbf{y}, \partial \mathcal{L} / \partial \mathbf{C}, \mathbf{g}_z)$$

$$\text{op3} = \text{TP}(\partial \mathcal{L} / \partial \mathbf{a}, \mathbf{y}, \mathbf{W})$$

$$\text{op4} = \text{backward}(\partial \mathcal{L} / \partial \mathbf{a}, \mathbf{y}, \mathbf{W}, \mathbf{g}_z)$$

$$\text{op5} = \text{backward}(\mathbf{x}, \partial \mathcal{L} / \partial \mathbf{b}, \mathbf{W}, \mathbf{g}_z)$$

$$\text{op6} = \text{TP}(\mathbf{x}, \partial \mathcal{L} / \partial \mathbf{b}, \mathbf{W})$$

$$\text{op7} = \text{TP}(\mathbf{x}, \mathbf{y}, \partial \mathcal{L} / \partial \mathbf{C})$$

 \implies

$$\partial \mathcal{L} / \partial \mathbf{x} = \text{op1}[1] + \text{op2}[1]$$

$$\partial \mathcal{L} / \partial \mathbf{y} = \text{op1}[2] + \text{op2}[2]$$

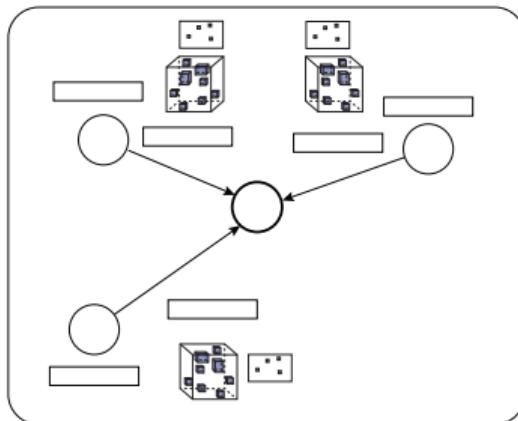
$$\partial \mathcal{L} / \partial \mathbf{W} = \text{op4}[3] + \text{op5}[3]$$

$$\partial \mathcal{L} / \partial \mathbf{g}_z = \text{op3} + \text{op6} + \text{op7}$$



Kernel Fusion

- Graph convolution has SpMM memory access pattern
- Fuse both kernels to save compute AND memory!



Algorithm Deterministic TP + Graph Convolution

Require: Graph $G = (V, E)$, $E[b] = (i_b, j_b)$

Require: Batch $x_1, \dots, x_{|V|}, y_1, \dots, y_{|E|}, W_1, \dots, W_{|E|}$

for segment_i \in schedule **do**

$(s, t) = E[k][0], E[k][1]$

Set $z_{\text{acc}} = 0$

▷ Parallel over Warps

for $b = 1 \dots |E|$ **do**

Execute segment subkernel sequence

$z_{\text{acc}} += z_{\text{smem}}$

if $b = |E|$ or $s < E[b+1][0]$ **then**

if s is first vertex processed by warp **then**

Send z_{acc} to fixup buffer.

else

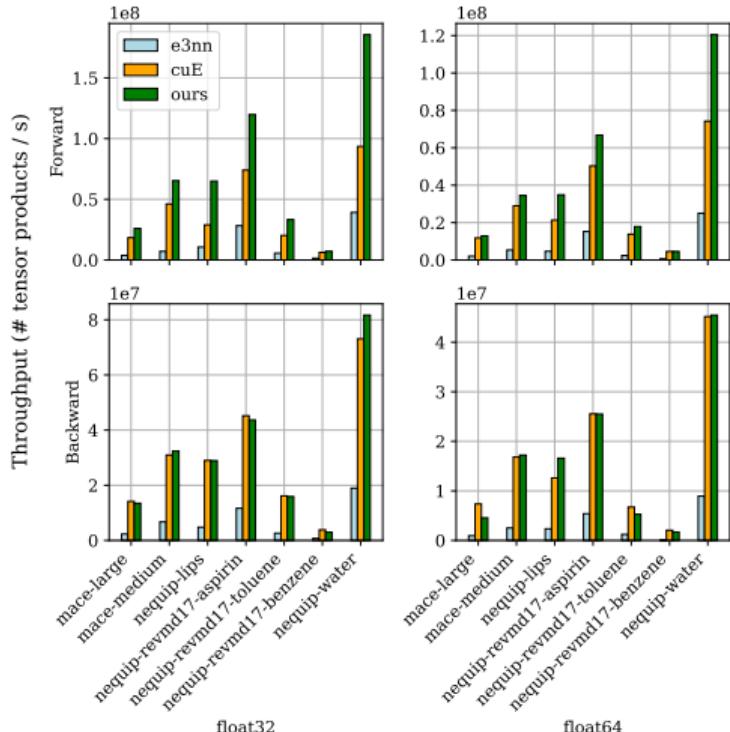
Store z_{acc} to global memory.

$z_{\text{acc}} = 0$

Execute fixup kernel.



Throughput (Unfused, A100)





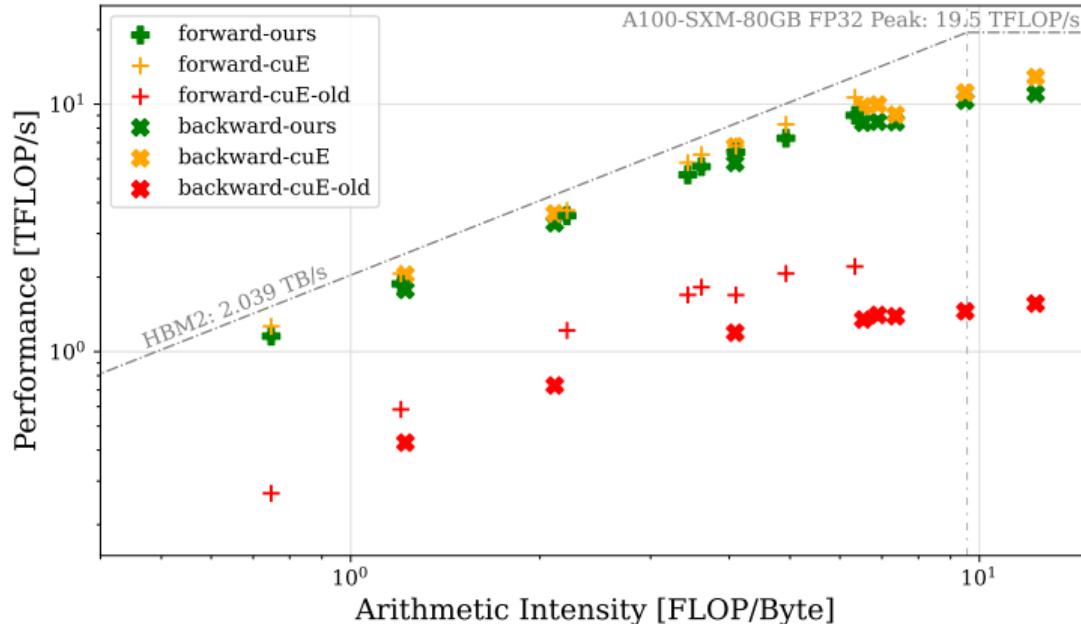
Cross-Platform Performance

GPU	forward			backward		
	e3nn	cuE	ours	e3nn	cuE	ours
A100	13	2.8	2.0	21	3.5	3.7
A5000	29	4.2	3.8	42	9.3	11
MI250x	41	-	3.0	128	-	15

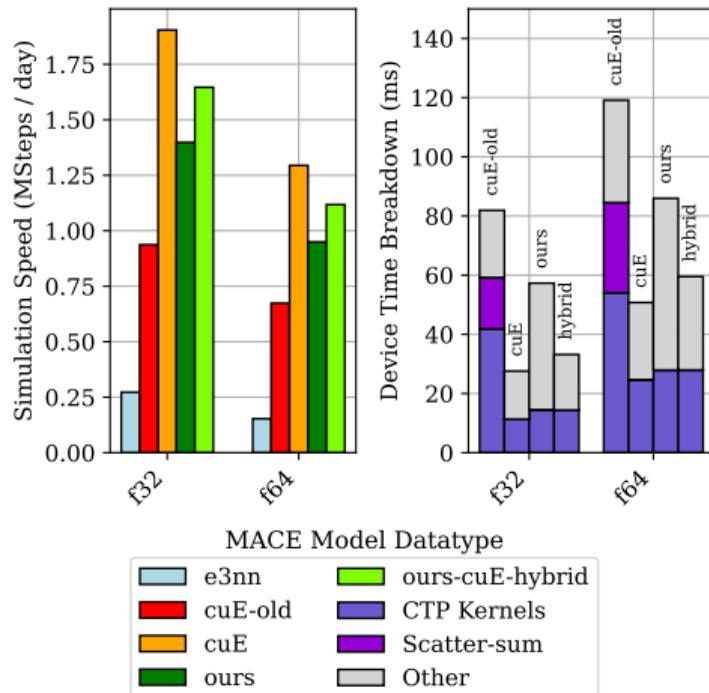
Table: MACE-large isolated tensor product runtime (ms), batch size 50K, FP32 unfused.



Roofline Analysis



Acceleration of MACE Foundation Model





Takeaways

- OpenEquivariance achieves SOTA performance, provides 5-6x end-to-end speedup for Nequip / MACE. 77 Github stars (and counting!)
- **Key Point:** Exploited structure in the tensor to engineer high performance kernels.
- **Related Work:** Graph neural networks have myriad applications; we study another useful kernel, SDDMM, in our IPDPS 2022 paper [BBD22].
- **Next:** a more general problem where we must **discover structure** in the sparse tensor.



Part 2: Randomized Sparse Tensor Decomposition

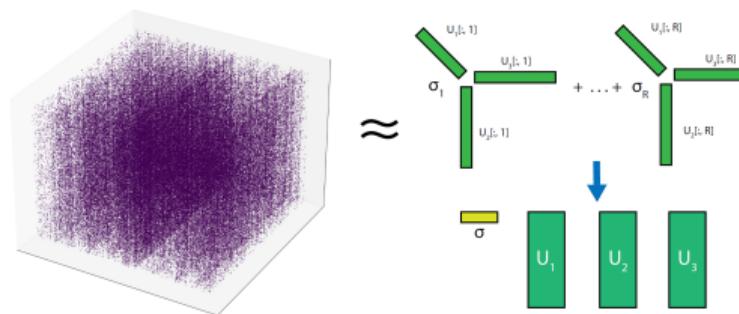


Motivating Application

- **Our goal:** efficiently solve an overdetermined linear least-squares problem

$$\min_{\mathbf{X}} \|\mathbf{AX} - \mathbf{B}\|_F$$

where $A = U_1 \odot \dots \odot U_N$ with $U_j \in \mathbb{R}^{I \times R}$.



- Key kernel in alternating least-squares Candecomp / PARAFAC (CP) decomposition [LK22].



ALS and the Multilinear Map

$$\min_{\mathbf{U}_j} \left\| \left[\bigodot_{k \neq j} \mathbf{U}_k \right] \cdot \mathbf{U}_j^\top - \text{mat}(\mathcal{T}, j)^\top \right\|_F$$

$$\min_{\mathbf{U}_2} \left\| \begin{array}{c} \mathbf{U}_3 \\ \odot \\ \mathbf{U}_1 \end{array} \cdot \begin{array}{c} \mathbf{U}_2^\top \\ - \\ \begin{array}{c} \cdot & \cdot \\ \cdot & \cdot \\ \cdot & \cdot \\ \cdot & \cdot \\ \cdot & \cdot \end{array} \\ \text{mat}(\mathcal{T}, 2) \end{array} \right\|_F \rightarrow \mathbf{U}_2 := \begin{array}{c} \mathbf{U}_2 \\ \vdots \\ \cdot \\ \cdot \\ \cdot \\ \cdot \end{array} \cdot \begin{array}{c} \mathbf{U}_3 \\ \odot \\ \mathbf{U}_1 \end{array} \cdot \mathbf{G}^+$$

MTTKRP



Randomized Linear Least-Squares

- **Sketch & Solve:** Apply short-wide sketching matrix S to both A and B , solve reduced problem

$$\min_{\tilde{\mathbf{X}}} \left\| S A \tilde{\mathbf{X}} - S B \right\|_F$$

- Want an (ε, δ) guarantee on solution quality: with high probability $(1 - \delta)$,

$$\left\| A \tilde{\mathbf{X}} - B \right\|_F \leq (1 + \varepsilon) \min_{\mathbf{X}} \| A \mathbf{X} - B \|$$

- In this talk: we will just **randomly sample** J rows of A and B to form a pair of much shorter matrices. Preserves structure in both input and output.



Our Contributions [Bha+23]

METHOD	SOURCE	ROUND COMPLEXITY (\tilde{O} NOTATION)
CP-ALS	[KB09]	$N(N + I)I^{N-1}R$
CP-ARLS-LEV	[LK22]	$N(R + I)R^N / (\varepsilon\delta)$
TNS-CP	[MAL22]	$N^3IR^3 / (\varepsilon\delta)$
GTNE	[MS22]	$N^2(N^{1.5}R^{3.5}/\varepsilon^3 + IR^2)/\varepsilon^2$
STS-CP	OURS	$N(NR^3 \log I + IR^2) / (\varepsilon\delta)$

- We build a data structure with runtime **logarithmic** in the height of the KRP and quadratic in R to sample from *leverage scores* of A .
- Yields the **STS-CP** algorithm: lower asymptotic runtime for randomized dense CP decomposition than recent SOTA methods (even more advantage for sparse tensors).



Intuition: Do I Need Every Row?

- Consider a univariate regression problem with 100,000 (x, y) points.
- This is a **highly-overdetermined** problem. Can pick a subset of points to perform fitting.

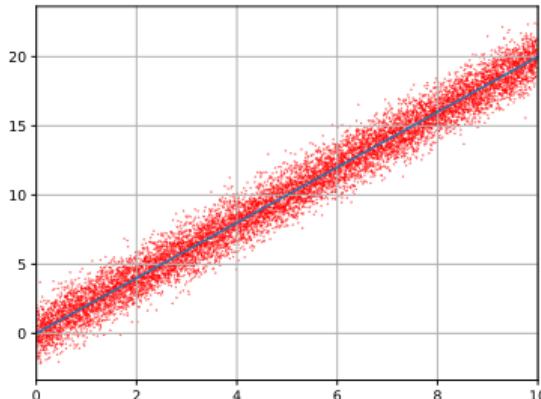


Figure: $y_i \sim x_i + \mathcal{N}(0, 1)$



Leverage Score Sampling

We will sample rows i.i.d. from \mathbf{A} according to the *leverage score distribution* on its rows. Given **reduced SVD** $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^\top$, define

$$\ell_i = \|\mathbf{U}[i, :]\|^2.$$

Theorem (Leverage Score Sampling Guarantees, [Mal22])

Suppose $\mathbf{S} \in \mathbb{R}^{J \times I}$ is a leverage-score sampling matrix for $\mathbf{A} \in \mathbb{R}^{I \times R}$, and define

$$\tilde{\mathbf{X}} := \arg \min_{\tilde{\mathbf{X}}} \left\| \mathbf{S} \mathbf{A} \tilde{\mathbf{X}} - \mathbf{S} \mathbf{B} \right\|_F$$

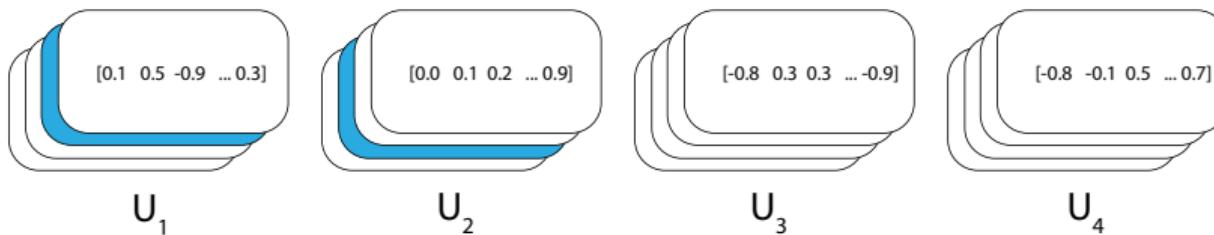
If $J \gtrsim R \max(\log(R/\delta), 1/(\varepsilon\delta))$, then with probability at least $1 - \delta$,

$$\left\| \mathbf{A} \tilde{\mathbf{X}} - \mathbf{B} \right\|_F \leq (1 + \varepsilon) \min_{\mathbf{X}} \left\| \mathbf{A} \mathbf{X} - \mathbf{B} \right\|_F.$$



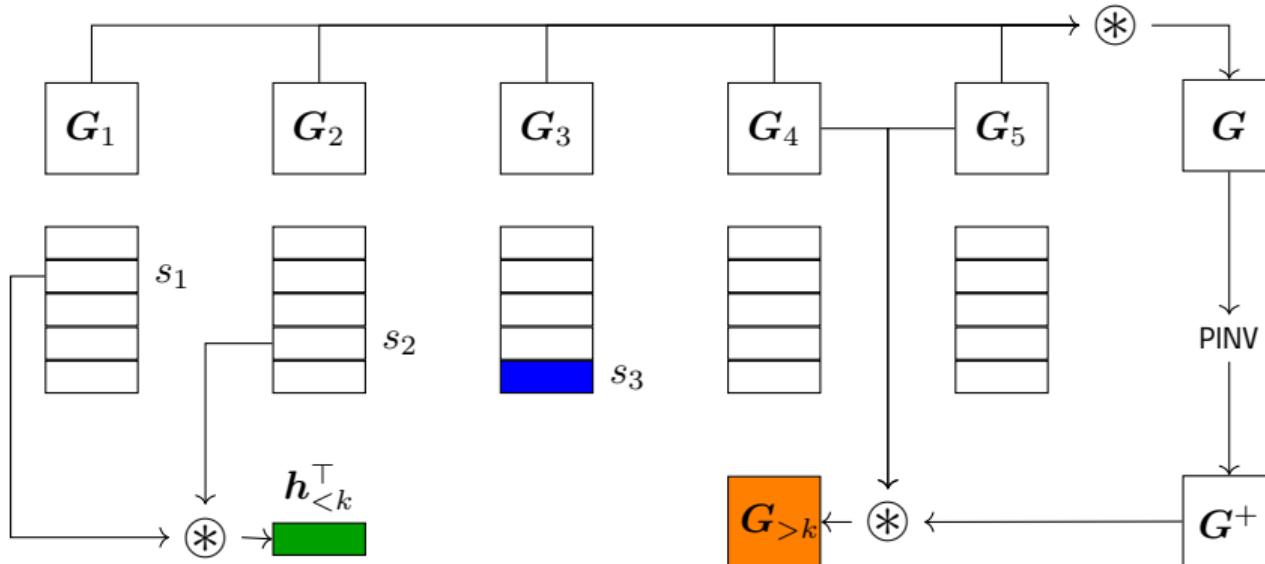
Implicit Leverage Score Sampling

- For $I = 10^7, N = 3$, matrix \mathbf{A} has 10^{21} rows. Can't even index rows with 64-bit integers. Instead: use identity $\ell_i = \mathbf{A}[i, :] (\mathbf{A}^\top \mathbf{A})^+ \mathbf{A}[i, :]^\top$.
- Draw a row from each of $\mathbf{U}_1, \dots, \mathbf{U}_N$, return Hadamard product.



- Let \hat{s}_j be a random variable for the row index drawn from \mathbf{U}_j . Assume $(\hat{s}_1, \dots, \hat{s}_N)$ jointly follows leverage score distribution on $\mathbf{U}_1 \odot \dots \odot \mathbf{U}_N$.

The Conditional Distribution of \hat{s}_k



$$p(\hat{s}_k = s_k \mid \hat{s}_{<k} = s_{<k}) \propto \langle \mathbf{h}_{<k} \mathbf{h}_{<k}^\top, \mathbf{U}_k [s_k, :]^\top \mathbf{U}_k [s_k, :], \mathbf{G}_{>k} \rangle$$



Key Sampling Primitive

$$\mathbf{q}[i] := C^{-1} \langle \mathbf{h}_{<k} \mathbf{h}_{<k}^\top, \mathbf{U}_k[i, :]^\top \mathbf{U}_k[i, :], \mathbf{G}_{>k} \rangle$$

- Imagine we magically had all entries of \mathbf{q} - how to sample? Initialize I bins, j -th has width $\mathbf{q}[j]$.
- Choose random real r in $[0, 1]$, find “containing bin” i such that

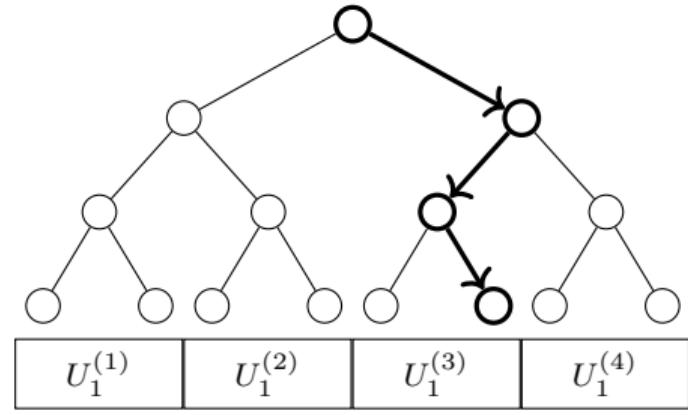
$$\sum_{j=0}^{i-1} \mathbf{q}[j] < r < \sum_{j=0}^i \mathbf{q}[j]$$



Binary Tree Inversion Sampling

- Locate bin via binary search (truncated to $\log(I/R)$ levels)
- Root: branch right iff $\sum_{j=0}^{I/2} \mathbf{q}[j] < r$
- Level 2: branch right iff

$$\sum_{j=0}^{I/2} \mathbf{q}[j] + \sum_{j=I/2}^{3I/4} \mathbf{q}[j] < r$$



Key: Can compute summations quickly if we cache information at each node!



Caching Partial Gram Matrices

Associate internal node v to an interval of rows $[S(v) \dots E(v)]$.

$$\sum_{j=S(v)}^{E(v)} \mathbf{q}[j] = \sum_{j=S(v)}^{E(v)} C^{-1} \langle \mathbf{h}_{<k} \mathbf{h}_{<k}^\top, \mathbf{U}_k[j, :]^\top \mathbf{U}_k[j, :], \mathbf{G}_{>k} \rangle$$

...

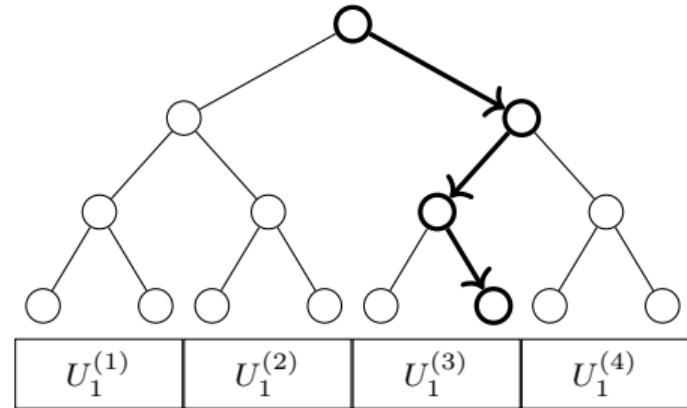
$$:= C^{-1} \langle \mathbf{h}_{<k} \mathbf{h}_{<k}^\top, \mathbf{G}^v, \mathbf{G}_{>k} \rangle$$

Can compute \mathbf{G}^v for ALL nodes v in time $O(IR^2)$, storage space $O(IR)$. Use BLAS-3 **syrk** calls in parallel to efficiently construct the tree.

Efficient Sampling after Caching



- At internal nodes, compute $C^{-1} \langle \mathbf{h}_{<k} \mathbf{h}_{<k}^\top, \mathbf{G}^v, \mathbf{G}_{>k} \rangle$ in $O(R^2)$ time (read normalization constant from root).
- At leaves, spend $O(R^3)$ time to compute remaining values of q . Can reduce to $O(R^2 \log R)$, see paper.
- Complexity per sample: $O(NR^2 \log I)$ (summed over all tensor modes).





Geometry Preservation

Define $D(S, A)$ of sketch S with respect to matrix A by

$$D(S, A) = \kappa(SQ) - 1$$

where Q is any orthonormal basis for the column space of A . Quantifies the distance preservation property of a sketch.

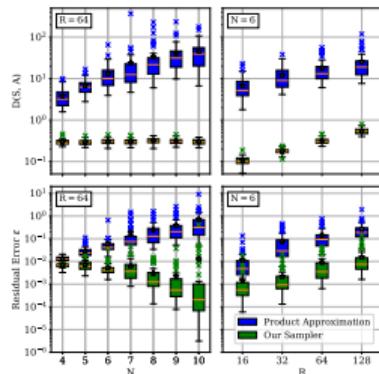


Figure: $D(S, A)$ as a function of KRP matrix count N and column count R , $J = 65, 536$. Green: our sampler. Blue: product approximation by [LK22].

Accuracy Comparison for Fixed Sample Count

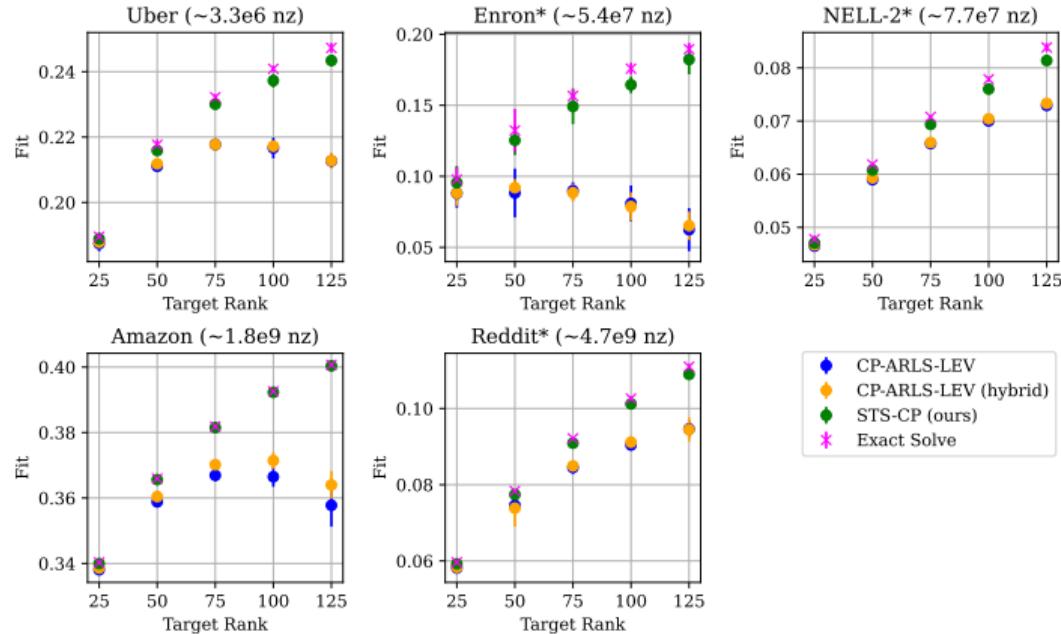


Figure: Sparse tensor ALS accuracy comparison for $J = 2^{16}$ samples, varied target ranks.

STS-CP Makes Faster Progress to Solution

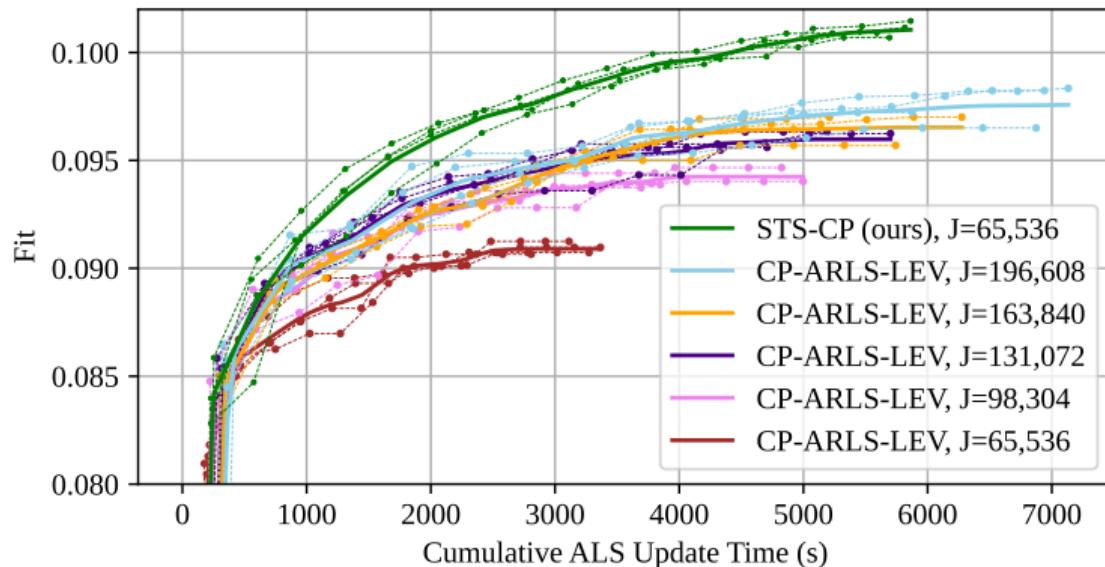


Figure: Fit vs. ALS update time, Reddit tensor, $R = 100$.



Takeaways

- Our data structure can downsample Khatri-Rao products with trillions of rows.
- Effectively preserves column-space geometry, offers OOM sample efficiency boost, 1.5-2.5x speedup over SOTA baseline methods.
- Data structures rely on common BLAS2, BLAS3 primitives, rendering them efficient.
- Compare to part 1: here, we exploit structure in the Khatri-Rao product, not the tensor.

Extension 1: Distributed Parallelism [Bha+24a]



- **SPAA 2024 Paper:** High-performance implementations of STS-CP and CP-ARLS-LEV in the distributed-memory parallel setting.
- Up to 11x speedup over SPLATT on hundreds of MPI ranks / thousands of CPU cores.
- We optimize both communication and computation in multiple stages of the solver.

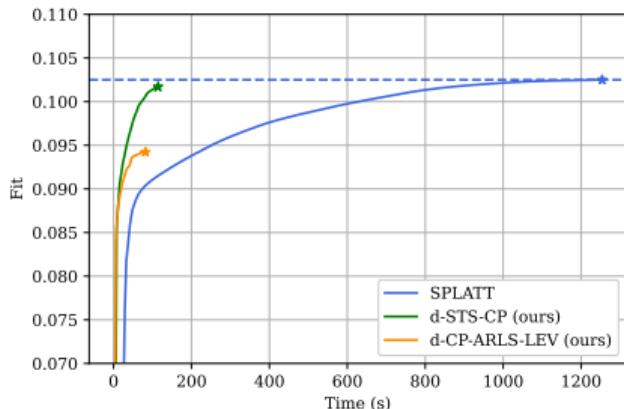
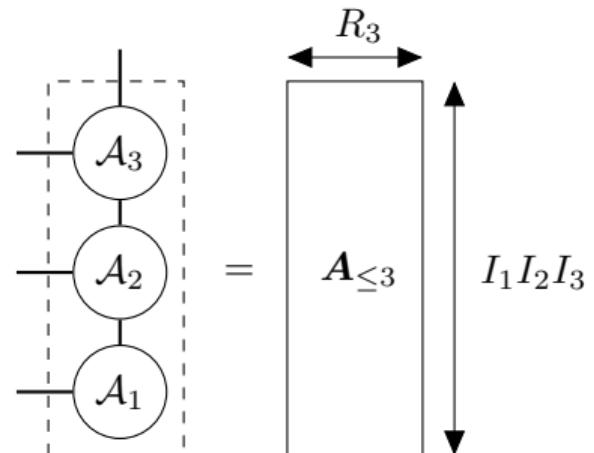


Figure: Accuracy vs. time, Reddit tensor, $R = 100$, 512 cores / 4 Perlmutter CPU nodes, 4.7B NNZ.

Extension 2: Sketched TT [Bha+24b]



- **NeurIPS 2024:** Apply theoretical machinery developed here to sketch *tensor trains* or *matrix product states*.
- Consider cores $\mathcal{A}_1, \dots, \mathcal{A}_j$, let $\mathbf{A}_{\leq j}$ be the *matricization* of the chain.
- When cores are in a special canonical form, exploit tree DS to sample $\mathbf{A}_{\leq j}$ with high efficiency.





Future Work: Inverse Problems

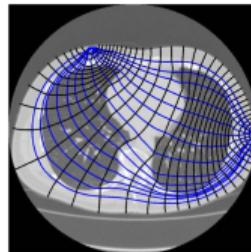


Figure: Electric field lines, human thorax (CC3.0 Wikimedia, Andy Adler).

- 2D EIT requires us to solve a linear least-squares problem of the form

$$\min_x \|(A_{11} \odot A_{12} + A_{21} \odot A_{22})x - b\|^2$$

where $A_{11}, A_{12}, A_{21}, A_{22}$ depend on geometry. Must adapt leverage-based sketch.

- More interesting problem: where to place electrodes along boundary to gain maximum information? *Sensor placement problem*, solve by volume (determinant) maximization.



Analysis and Future Directions



The Upshot

- This presentation focused on a single operation... but we used techniques from
 - ILP-maximizing, fine-grained GPU kernel engineering
 - Randomized algorithm theory
 - Distributed-memory communication analysis
- We explored parallel computing motifs that include
 - Dense matrix multiply (warp-level)
 - Multiple sparse matrix / tensor primitives
 - Particle simulation / molecular dynamics
- We accelerated applications ranging from chemistry foundation models to multidimensional pattern mining / data analysis.
- A lot of mileage for one simple primitive!

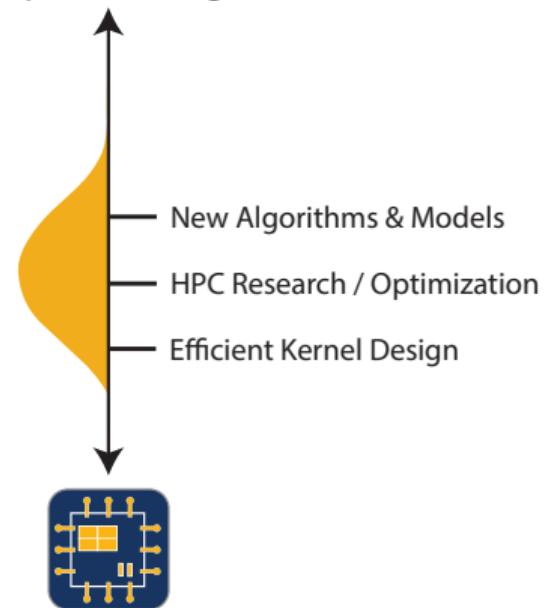


Future Work

- **Up the Stack:** How do we co-design new MLIP models / algorithms to exploit accelerated primitives?
- **Down the Stack:** How do we design systems (e.g. ML compilers) to automate kernel engineering while delivering high performance?
- **Scientific / Industry Collaborations:** So much needs acceleration - got a problem? Let's talk.



P / NP, PSPACE...
 $O(\exp(N))$, $O(N \log N)$...





Acknowledgements - Thank you to...



Everyone in Attendance!

(and your taxpayer dollars funding this research)



<https://www.xkcd.com/1403/>



My Advisers



Aydin Buluç



James Demmel

Committee, Collaborators, and Fellow Grads



Committee Members: Michael Lindsey and Katherine Yelick– thank you!

Collaborators: Russell Castro, Austin Glover, Laura Grigori, Riley Murray, Guillaume Rabusseau and Beheshteh Rakhshan, with very special thanks to Osman Malik.

Members of PASSION and BeBOP: Benjamin Brock, Alok Tripathy, Giulia Guidi, Roger Hsiao, Tianyu Liang, Gabriel Raulet, Yen-Hsiang Chang.



Supporters Over Many Years

Undergraduate / internship advisers: Chris Umans, Rose Yu, Mikhail Shapiro, and Adam Sheffer while I was at Caltech, and Federico Busato at NVIDIA!

Members of the CSGF Program: Koby Hayashi, Caleb Ju, Grace Wei, Mansi Sakarvadia, Santiago Vargas, Albert Muesalian, Gabby Jones, Julian Bellavita, Kiran Eiden, Mary Francis LaPorte, Olorundamilola Kazeem, Ethan Epperly, Katherine Keegan, and so many others.

Especially: Allison Wang, Anant Kale, Hanwen Zhang, Emily Pan, Daniel Mark. And of course, Aditi Lahiri. Thank you for everything.



My Family





La fin!



References I

- [Batz+22] Simon Batzner, Albert Musaelian, Lixin Sun, Mario Geiger, Jonathan P. Mailoa, Mordechai Kornbluth, Nicola Molinari, Tess E. Smidt, and Boris Kozinsky. “E(3)-equivariant graph neural networks for data-efficient and accurate interatomic potentials”. In: *Nature Communications* 13.1 (May 2022), p. 2453. ISSN: 2041-1723. doi: 10.1038/s41467-022-29939-5. URL: <https://doi.org/10.1038/s41467-022-29939-5>.
- [Bat+24] Illyes Batatia et al. *A foundation model for atomistic materials chemistry*. 2024. arXiv: 2401.00096 [physics.chem-ph]. URL: <https://arxiv.org/abs/2401.00096>.
- [BBD22] Vivek Bharadwaj, Aydin Buluç, and James Demmel. “Distributed-Memory Sparse Kernels for Machine Learning”. In: *2022 IEEE International Parallel and Distributed Processing Symposium*. Los Alamitos, CA, USA: IEEE Computer Society, June 2022, pp. 47–58. doi: 10.1109/IPDPS53621.2022.00014.
- [Bha+23] Vivek Bharadwaj, Osman Asif Malik, Riley Murray, Laura Grigori, Aydin Buluç, and James Demmel. “Fast Exact Leverage Score Sampling from Khatri-Rao Products with Applications to Tensor Decomposition”. In: *Thirty-seventh Conference on Neural Information Processing Systems*. Dec. 2023.



References II

- [Bha+24a] Vivek Bharadwaj, Osman Asif Malik, Riley Murray, Aydin Buluç, and James Demmel. “Distributed-Memory Randomized Algorithms for Sparse Tensor CP Decomposition”. In: *Proceedings of the 36th ACM Symposium on Parallelism in Algorithms and Architectures*. SPAA ’24. Nantes, France: Association for Computing Machinery, May 2024, pp. 155–168. ISBN: 9798400704161. doi: 10.1145/3626183.3659980. URL: <https://doi.org/10.1145/3626183.3659980>.
- [Bha+24b] Vivek Bharadwaj, Beheshteh T. Rakhshan, Osman Asif Malik, and Guillaume Rabusseau. “Efficient Leverage Score Sampling for Tensor Train Decomposition”. In: *Thirty-eighth Conference on Neural Information Processing Systems*. Dec. 2024. arXiv: 2406.02749 [cs.DS].
- [Bha+25] Vivek Bharadwaj, Austin Glover, Aydin Buluç, and James Demmel. “An Efficient Sparse Kernel Generator for $O(3)$ -Equivariant Deep Networks”. In: *Proceedings of the 3rd SIAM Conference on Applied and Discrete Computational Algorithms*. July 2025. URL: <https://arxiv.org/abs/2501.13986>.
- [Che+20] Ke Chen, Qin Li, Kit Newton, and Stephen J. Wright. “Structured Random Sketching for PDE Inverse Problems”. In: *SIAM Journal on Matrix Analysis and Applications* 41.4 (2020), pp. 1742–1770. doi: 10.1137/20M1310497.



References III

- [Cor+23] Gabriele Corso, Hannes Stärk, Bowen Jing, Regina Barzilay, and Tommi S. Jaakkola. "DiffDock: Diffusion Steps, Twists, and Turns for Molecular Docking". In: *The Eleventh International Conference on Learning Representations*. 2023.
- [Dao+22] Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. "FlashAttention: Fast and Memory-Efficient Exact Attention with IO-Awareness". In: *Advances in Neural Information Processing Systems*. Ed. by S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh. Vol. 35. Curran Associates, Inc., 2022, pp. 16344–16359.
- [KB09] Tamara G. Kolda and Brett W. Bader. "Tensor Decompositions and Applications". In: *SIAM Review* 51.3 (Aug. 2009). Publisher: Society for Industrial and Applied Mathematics, pp. 455–500. ISSN: 0036-1445. DOI: 10.1137/07070111X.
- [LK22] Brett W. Larsen and Tamara G. Kolda. "Practical Leverage-Based Sampling for Low-Rank Tensor Decomposition". In: *SIAM Journal on Matrix Analysis and Applications* 43.3 (2022), pp. 1488–1517.



References IV

- [Mal22] Osman Asif Malik. "More Efficient Sampling for Tensor Decomposition With Worst-Case Guarantees". In: *Proceedings of the 39th International Conference on Machine Learning*. Vol. 162. Proceedings of Machine Learning Research. PMLR, July 2022, pp. 14887–14917.
- [MS22] Linjian Ma and Edgar Solomonik. "Cost-efficient Gaussian tensor network embeddings for tensor-structured inputs". In: *Advances in Neural Information Processing Systems*. Vol. 35. Curran Associates, Inc., 2022, pp. 38980–38993.
- [Smi+17] Shaden Smith, Jee W. Choi, Jiajia Li, Richard Vuduc, Jongsoo Park, Xing Liu, and George Karypis. *FROSTT: The Formidable Repository of Open Sparse Tensors and Tools*. 2017. URL: <http://frostdt.io/>.