



BERKELEY LAB

Bringing Science Solutions to the World



U.S. DEPARTMENT OF
ENERGY

Office of Science

Engineering Fast Kernels for Rotation-Equivariant Chemical Foundation Models

Vivek Bharadwaj, DOE CSGF Program Review 2025



Sparse Tensors, Multilinear Maps, and Me

About me: computer science PhD student at UC Berkeley.

Interests: sparse linear algebra, multilinear maps and tensor kernels to accelerate scientific / ML workloads.

Some past projects:

- Algorithms for high-dimensional “tensor SVD”.
- Scaling bulk-synchronous linear algebra on large clusters.
- Distributed-memory graph algorithms.

$$\begin{bmatrix} \mathbf{z}_{:1} \\ \mathbf{z}_{:2} \\ \mathbf{z}_{:3} \end{bmatrix} = \mathbf{A} \cdot \begin{bmatrix} \mathbf{x}_{:1} \\ \mathbf{x}_{:2} \\ \mathbf{x}_{:3} \end{bmatrix}$$
$$\mathbf{z}_{:i} = \mathbf{A} \cdot (\mathbf{x}_{:i} \otimes \mathbf{y}_{:i})$$
$$\mathbf{z}_{ji} = \sum_{u,v} \mathcal{A}_{juv} \mathbf{x}_{ui} \mathbf{y}_{vi}$$

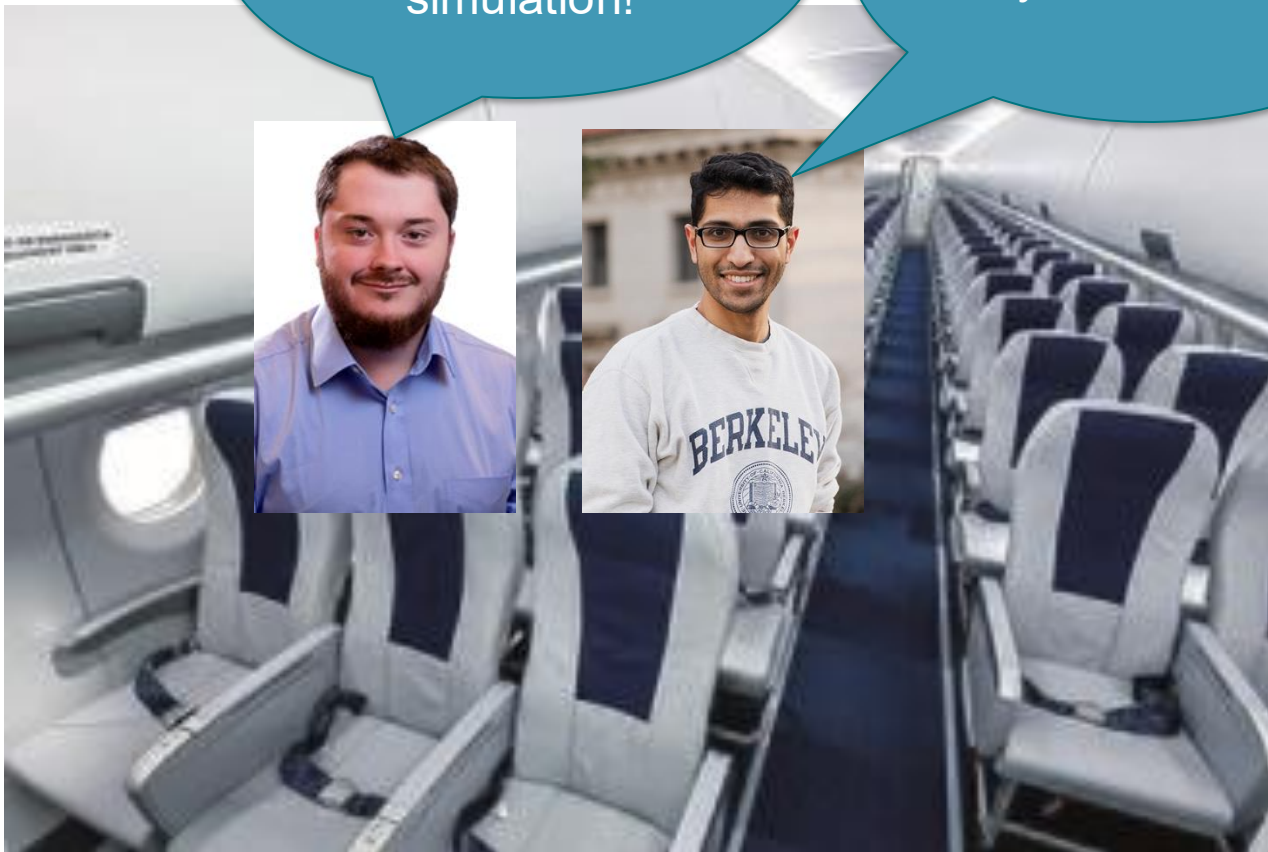
$$\begin{bmatrix} \mathbf{z}_{:1} \\ \mathbf{z}_{:2} \\ \mathbf{z}_{:3} \end{bmatrix} = \mathcal{A} \cdot \begin{bmatrix} \mathbf{x}_{:1} \\ \mathbf{x}_{:2} \\ \mathbf{x}_{:3} \end{bmatrix} \cdot \begin{bmatrix} \mathbf{y}_{:1} \\ \mathbf{y}_{:2} \\ \mathbf{y}_{:3} \end{bmatrix}$$

 **AM NOT:** A chemist, computational or otherwise. 

My Flight to the 2024 CSGF Program Review

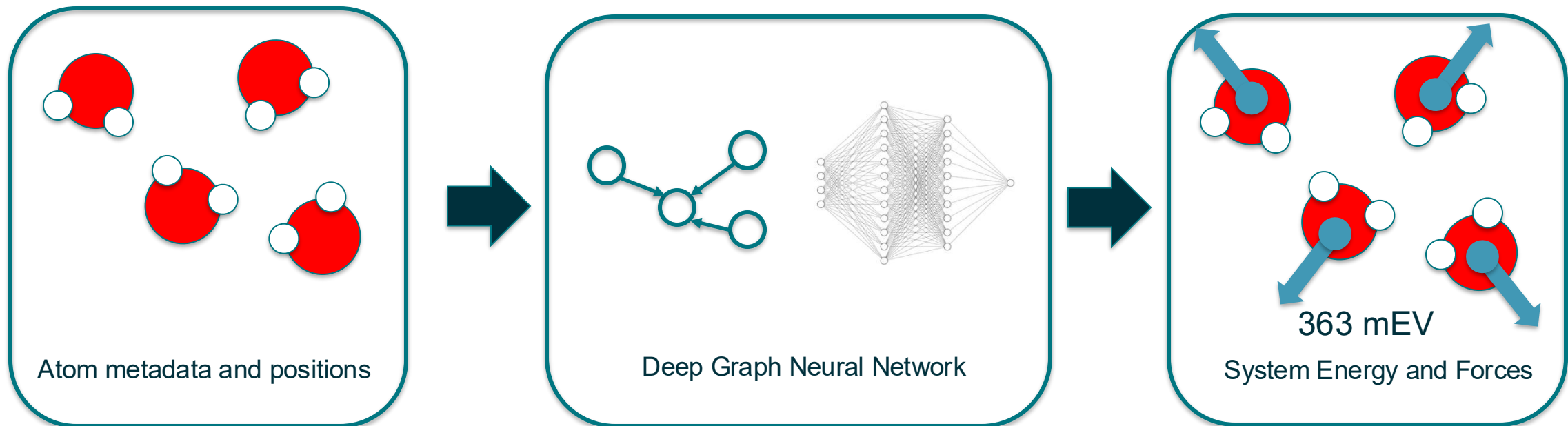
We study MLIAPs
for chemical
simulation!

so which atom is
your favorite



Tristan Maxson (CSGF Y3) generously spent hours telling me about *machine learning interatomic potentials*.

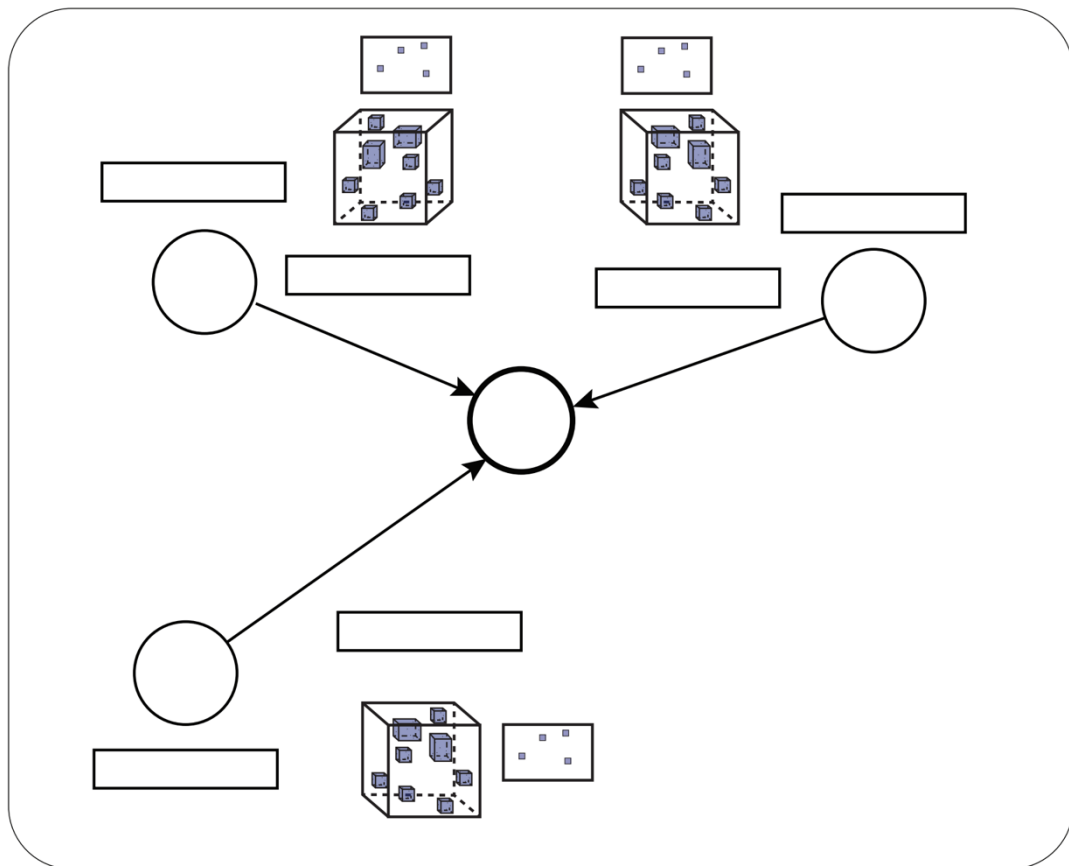
Interatomic Potentials via Graph Neural Networks



Goal: Use a message-passing graph neural network to predict atomic energies and forces.

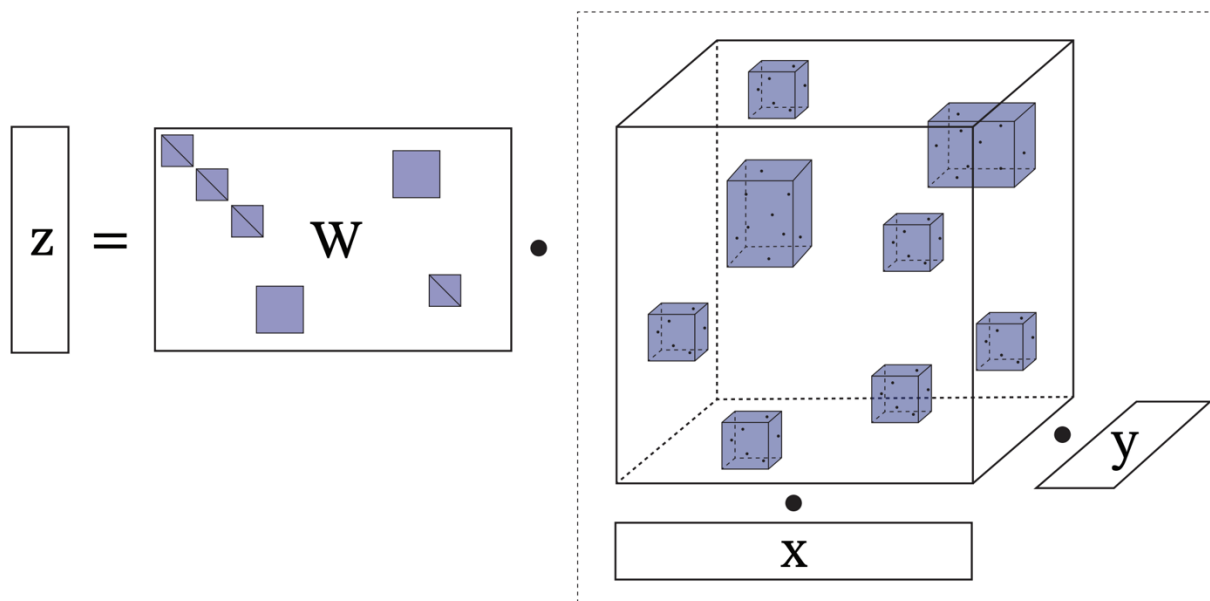
Constraint: Want our network to respect certain physical laws (e.g., symmetry, energy conservation).

Respecting Physical Symmetries



- Message-passing GNNs generate feature vectors for each node and edge.
- **$O(3)$ -Equivariance:** If 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



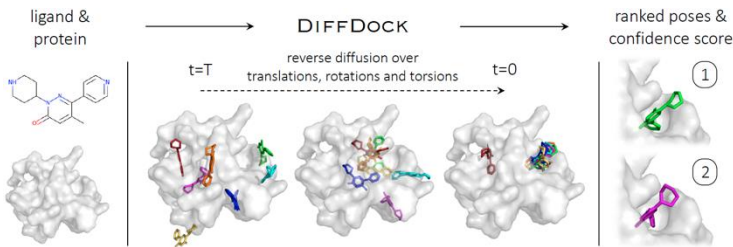
Structure of the sparse tensor is known when GNN is initialized.

Typically, need to execute **millions** of these operations.

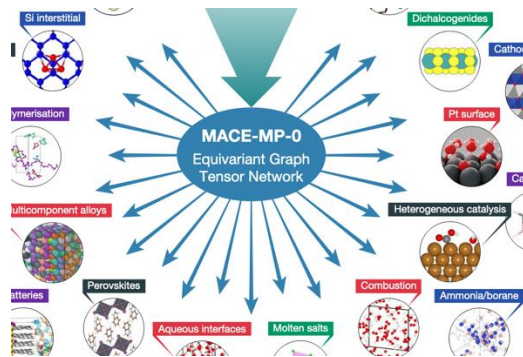
Problem: GPUs are optimized for dense GEMM!

CG tensor product can run at <10% GPU utilization, consumes >80% of model runtime.

Some Large Equivariant Neural Network Projects



Bowen Jing (Y4) & others



Samuel Blau (alum '16) & others

We simulated
an HIV capsid
on 5120 GPUs.



Albert Musaelian (alum'23) and
colleagues nominated for the
2023 **Gordon Bell prize**.

OpenEquivariance: Turbocharging CG Performance

- We present OpenEquivariance: an open-source GPU kernel generator for the CG tensor product
- **OOM speedup over best open-source codes**, on-par w/ NVIDIA's **closed-source** package (but 2x speedup against their earlier version!)
- 5-6x speedup for MACE / Nequip, 10x+ smaller memory footprint. 76 Github stars and counting!



Me*



Austin
Glover*

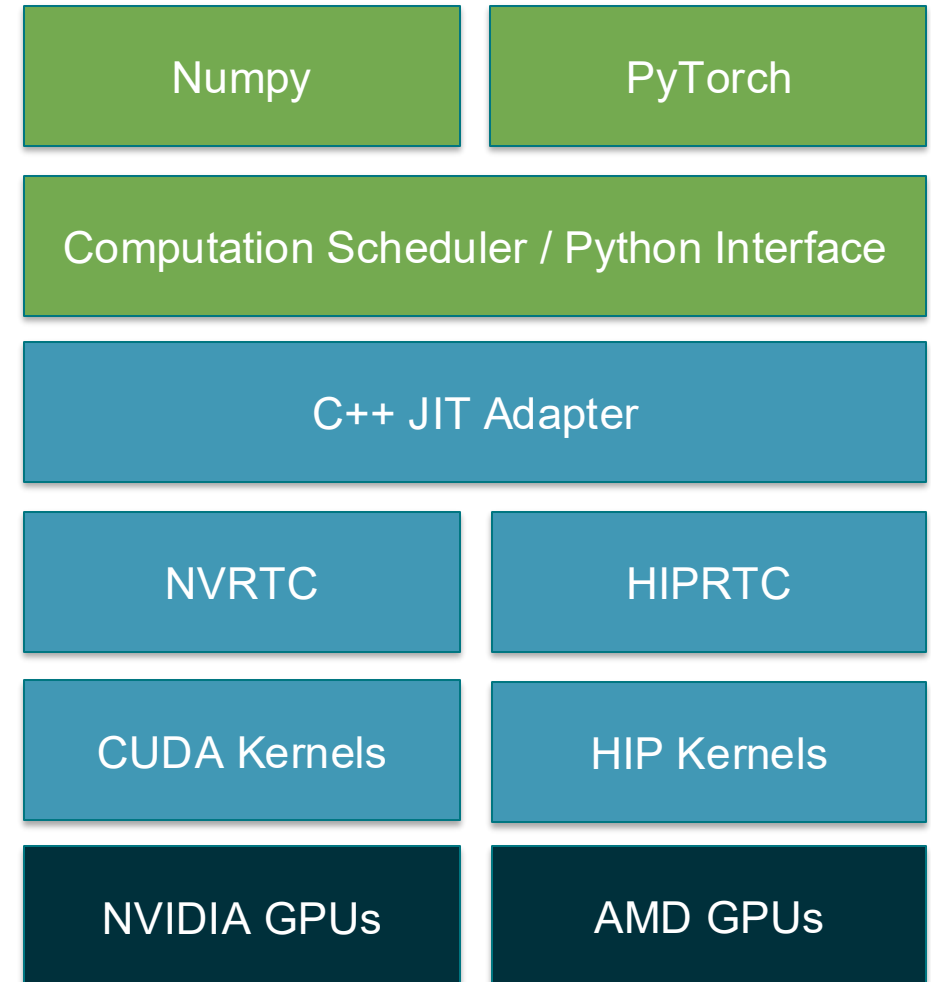


Aydın
Buluç



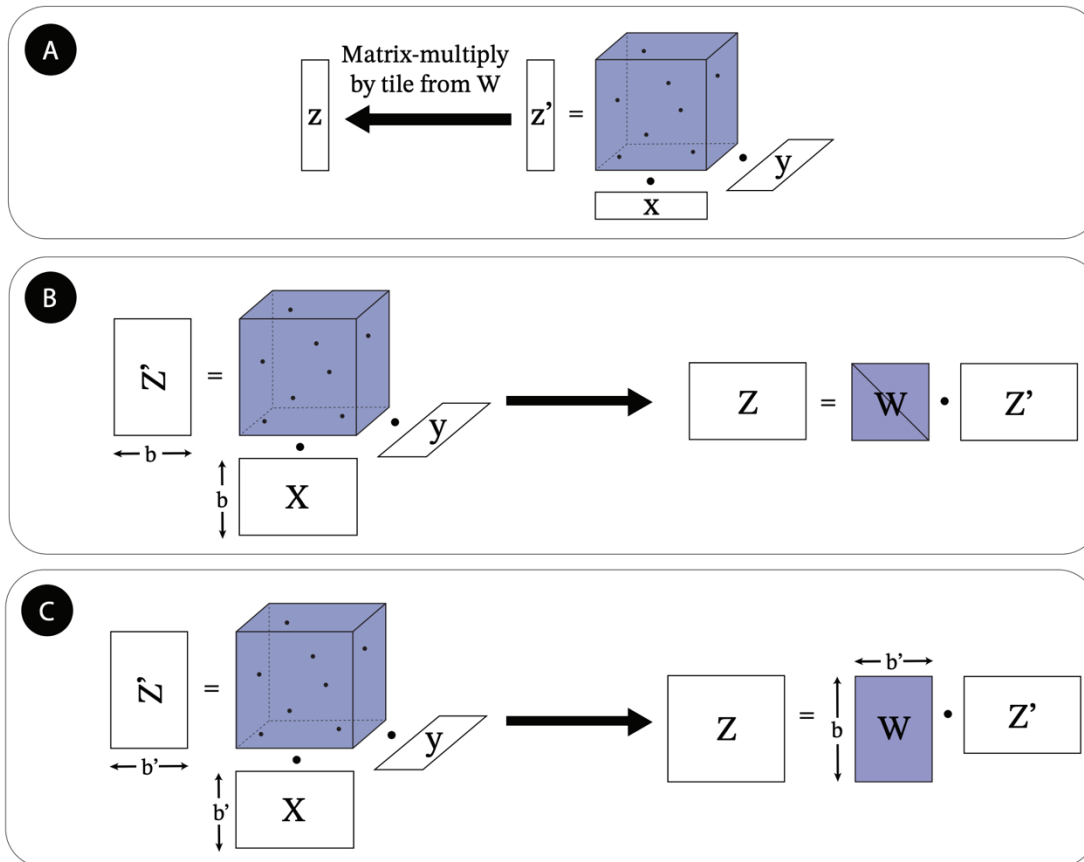
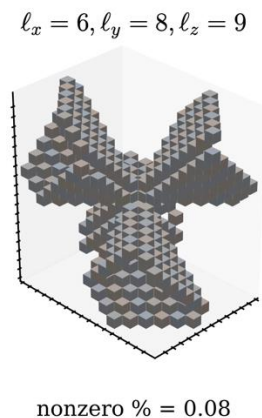
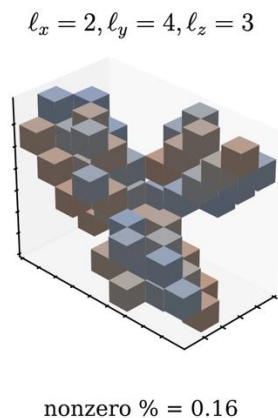
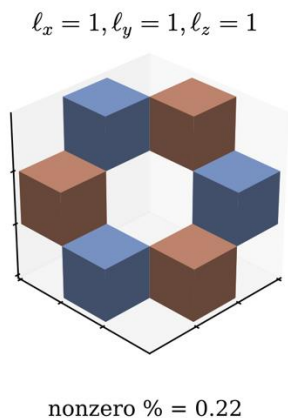
James
Demmel

*=equal effort



Subkernels of the CG Tensor Product

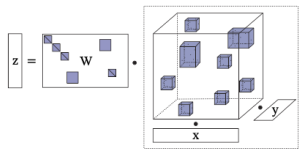
- Break the computation into segments, one for each nonzero block.
- Each segment contracts x , y with a block, matrix-multiplies result by tile from W .



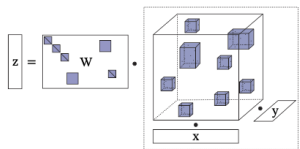
A Roadmap to Efficient CG Kernels

Assign each batch element to distinct GPU warp

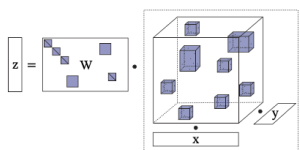
Warp 1



Warp 2



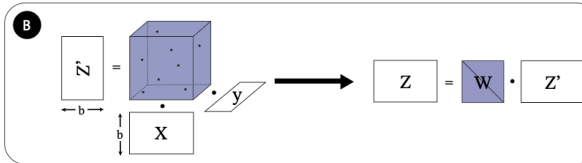
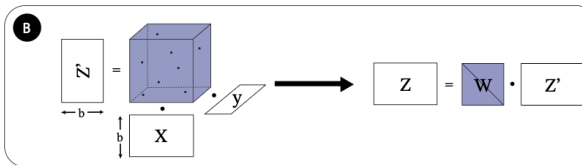
Warp 3



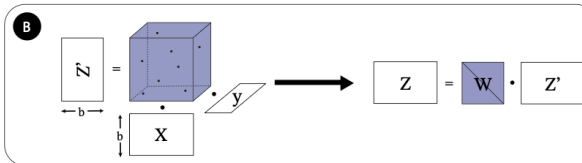
⋮

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] += 0x1.6a09e6000000p-2f * x[0] * y[0];
z[0] += -0x1.6a09e6000000p-2f * x[1] * y[0];
z[2] += 0x1.d363d0000000p-2f * x[1] * y[0];
z[1] += -0x1.d363d0000000p-2f * x[2] * y[0];
z[4] += -0x1.6a09e6000000p-1f * x[3] * y[0];
z[3] += 0x1.6a09e6000000p-1f * x[4] * y[0];
z[5] += -0x1.d363d0000000p-2f * x[4] * y[0];
z[4] += 0x1.d363d0000000p-2f * x[5] * y[0];
z[6] += -0x1.6a09e6000000p-2f * x[5] * y[0];
z[5] += 0x1.279a74000000p-1f * x[1] * y[1];
z[4] += 0x1.279a74000000p-2f * x[2] * y[1];
z[2] += -0x1.279a74000000p-2f * x[4] * y[1];
z[1] += -0x1.279a74000000p-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, 0f3F08D677;
mul.f32 %f4187, %f4186, %f4096;
fma.rn.f32 %f4188, %f4187, %f4089, %f4146;
fma.rn.f32 %f4189, %f4100, %f4089, %f4188;
mul.f32 %f4190, %f4133, 0f3EA79762;
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;
```

Forward Pass Warp-Level 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 $(i, j, k, v) \in \text{nz}(\mathcal{P})$ **do** ▷ Unroll via JIT
 $\mathbf{z}_{\text{reg}}[k] \mathrel{+}= v \cdot \mathbf{x}_{\text{reg}}[i] \cdot \mathbf{y}_{\text{reg}}[j]$

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

Backward calculation is similar, but requires three update equations instead of one.

Optimizing Model Force Predictions

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

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

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

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

Novel Identities for Second Partial Derivatives

- PyTorch is happy to calculate a 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.

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

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

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

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

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

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

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



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

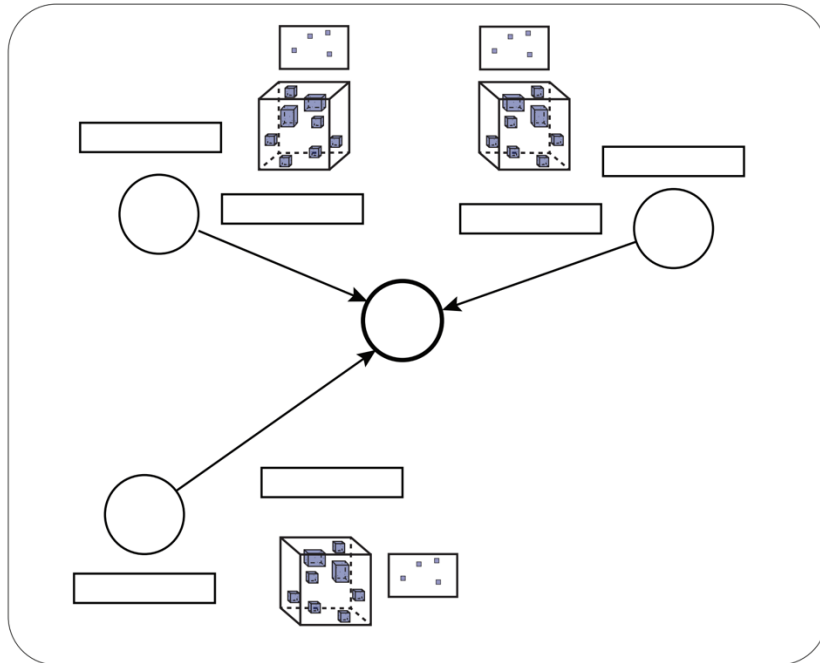
$\partial\mathcal{L}/\partial\mathbf{y} = \text{op1} [2] + \text{op2} [2]$

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

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

Graph Convolution with Kernel Fusion

- Graph convolution has a well-studied SpMM memory access pattern.
- Fuse both CG tensor product and graph convolution to save compute AND memory!



Algorithm Deterministic TP + Graph Convolution

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

Require: Batch $\mathbf{x}_1, \dots, \mathbf{x}_{|V|}$, $\mathbf{y}_1, \dots, \mathbf{y}_{|E|}$, $\mathbf{W}_1, \dots, \mathbf{W}_{|E|}$

for segment _{i} \in schedule **do**

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

Set $\mathbf{z}_{\text{acc}} = 0$

▷ Parallel over Warps

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

Execute segment subkernel sequence

$\mathbf{z}_{\text{acc}} \mathrel{+}= \mathbf{z}_{\text{smem}}$

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

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

Send \mathbf{z}_{acc} to fixup buffer.

else

Store \mathbf{z}_{acc} to global memory.

$\mathbf{z}_{\text{acc}} = 0$

Execute fixup kernel.

Experiments on DOE Flagship Systems

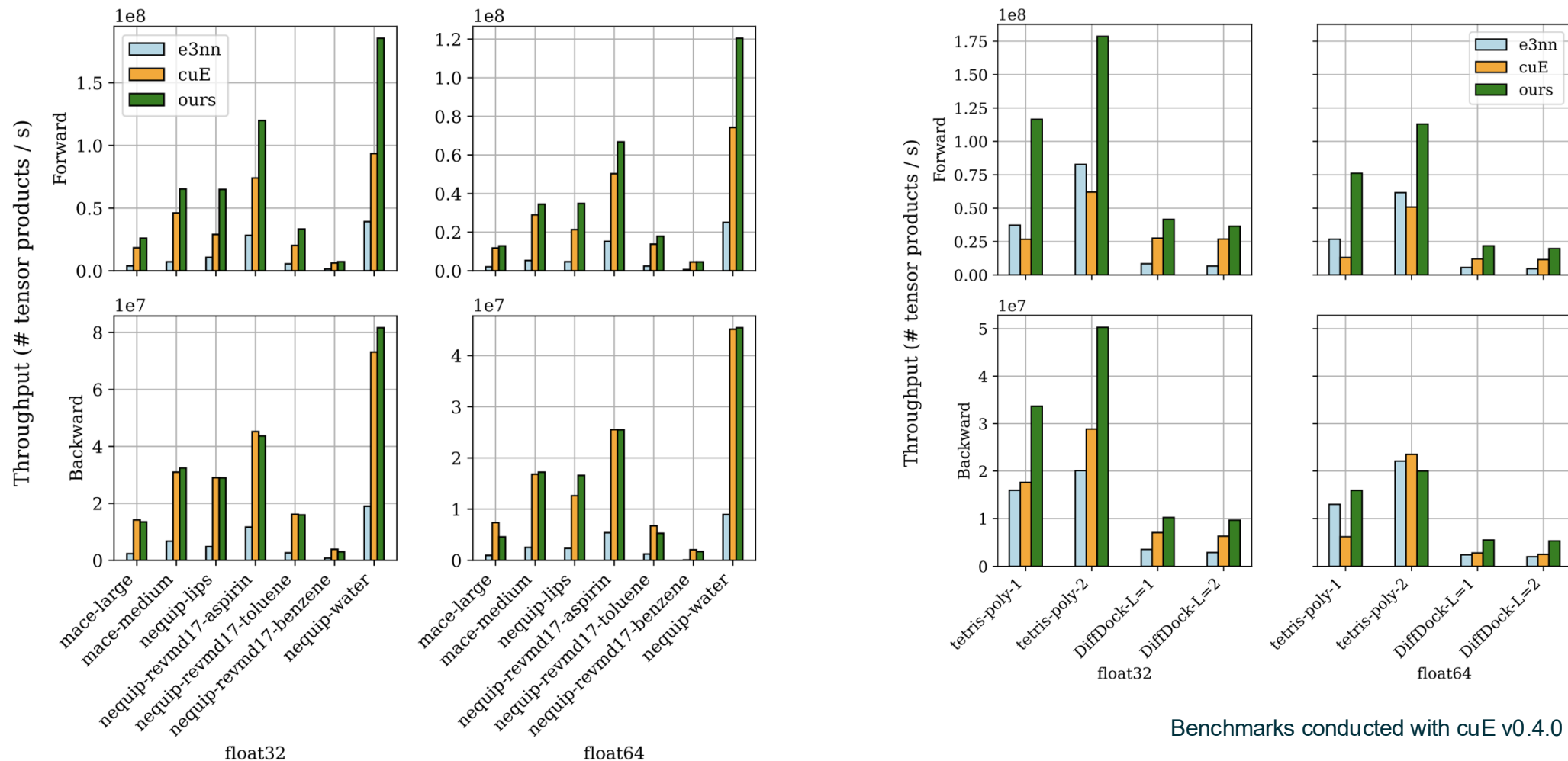


NERSC Perlmutter, A100 GPUs



OLCF Frontier, MI250x GPUs

A100 Tensor Product Throughput

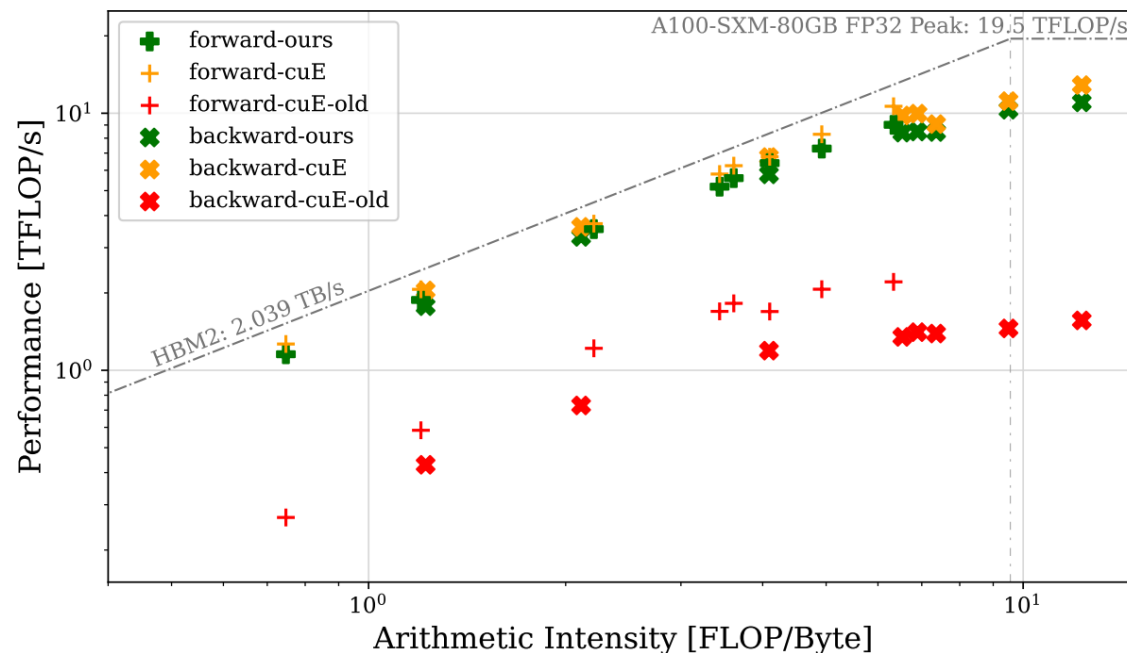


Benchmarks conducted with cuE v0.4.0

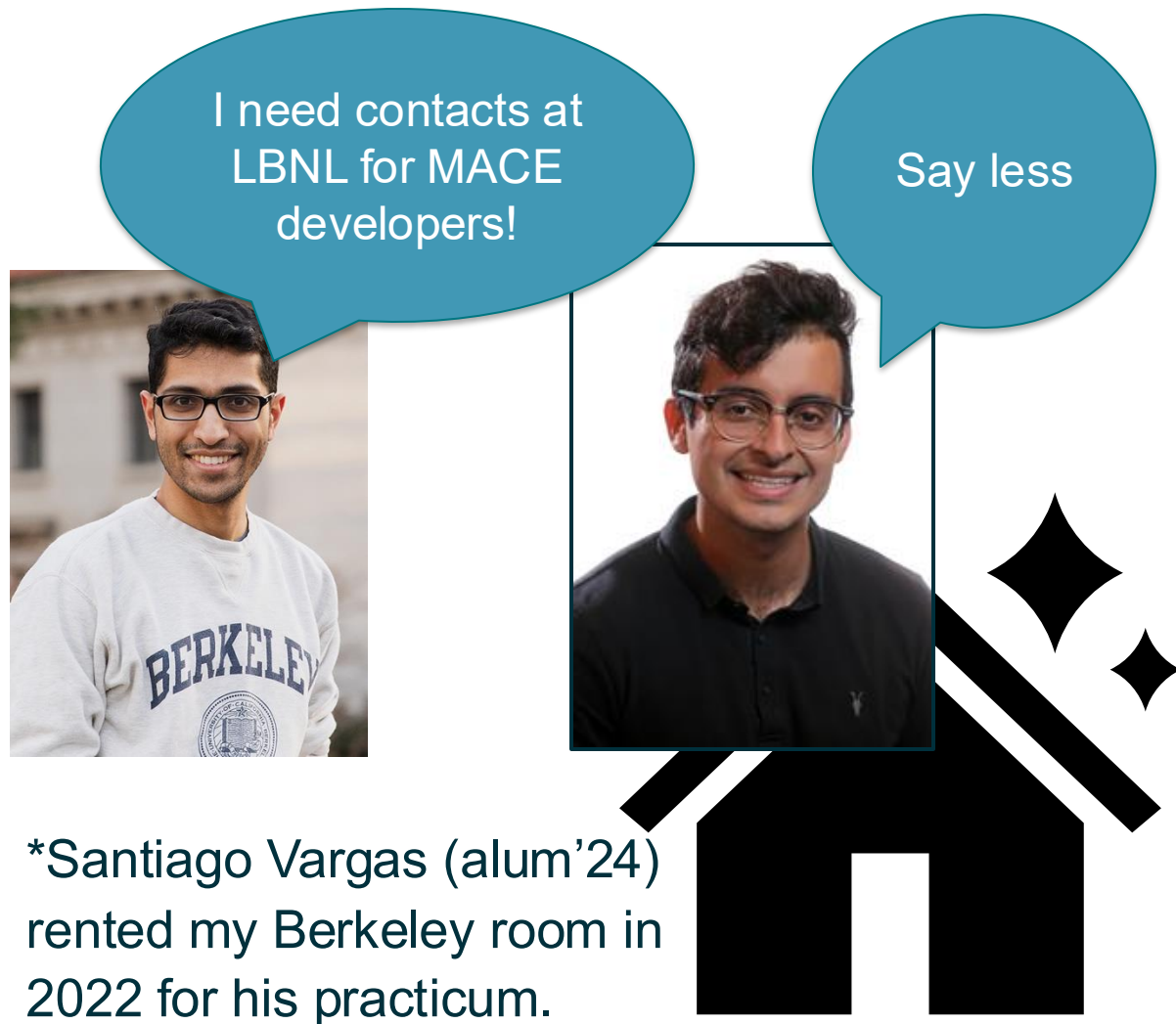
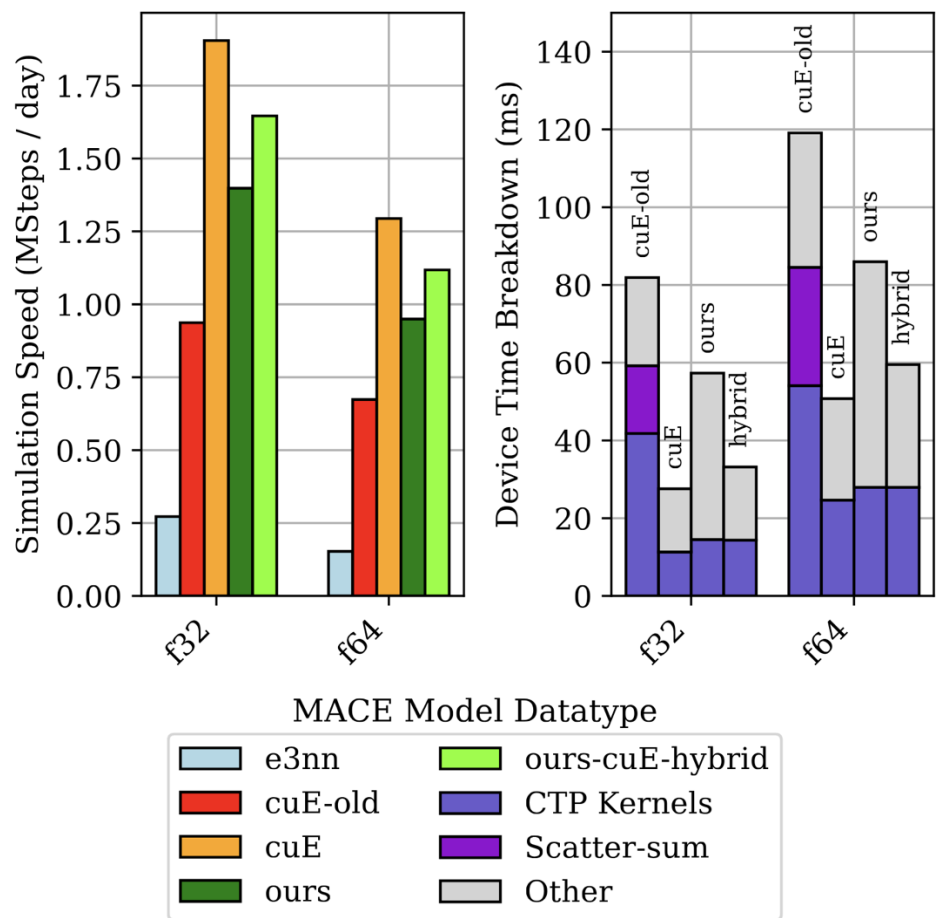
Cross-Platform Performance and Roofline Analysis

| 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 |

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



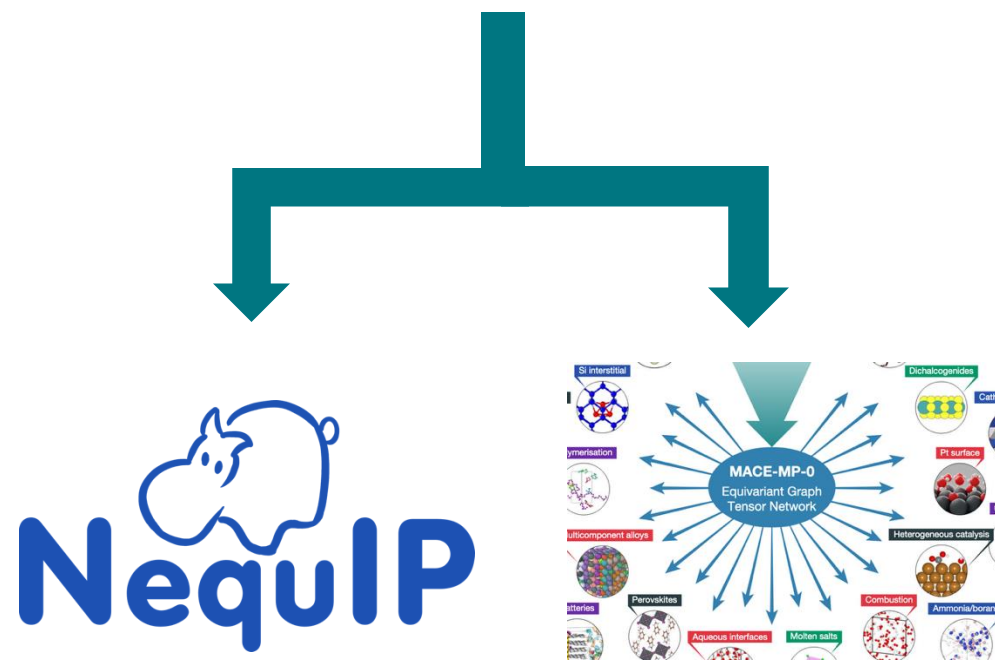
Accelerating MACE (Lawrence Berkeley Lab)



Conclusions and Further Work

- OpenEquivariance is now **officially integrated** into Nequip and MACE.
- 5-6x E2E speedup, OOM memory reduction via:
 - Synchronization-light parallelization strategies.
 - Careful management of GPU shared memory.
 - Register caching / loop unrolling.
 - Kernel fusion, and more!
- Tentative future work: numerical precision tuning, JAX wrappers, AMD performance tuning.

```
pip install openequivariance
```



The Real Atoms were the Friends Along the Way ☐

(that doesn't scan) Thank you Lindsey, Kris, Michelle, Robyn, and the entire Krell staff!



Questions?

Scan the code for the package and paper.

