

An Efficient Sparse Kernel Generator for $O(3)$ -Equivariant Deep Networks

Vivek Bharadwaj, Austin Glover, Aydin Buluç, James Demmel

Presented by Nicole Hao

December 14, 2025

Overview

① Background

② Solution

③ Experimental Results & Performance Gains

④ Key Takeaways & Future Directions

Geometric Deep Learning in Atomic Simulation

- Input: Atomic metadata and positions.
- Process: A message-passing graph neural network (MP-GNN).
- Output: 1. System energy (in milli-electronvolts, meV) 2. Atomic forces.

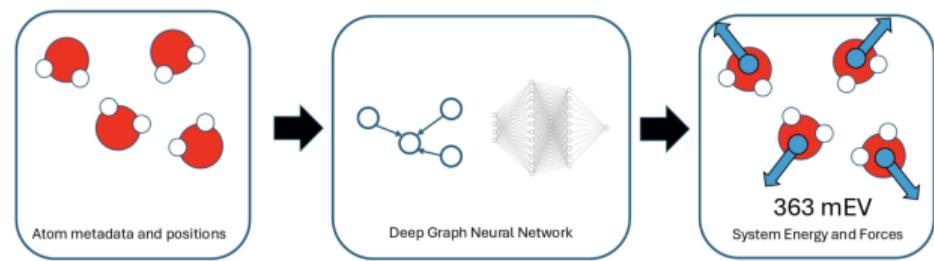


Figure: Interatomic Potential Calculation

Why do we need geometric deep learning in molecular simulations?

Geometric deep learning achieves state-of-the-art performance for fast interatomic potential calculation. A key primitive is the **message-passing graph neural network**.

Message-passing Graph Neural Network

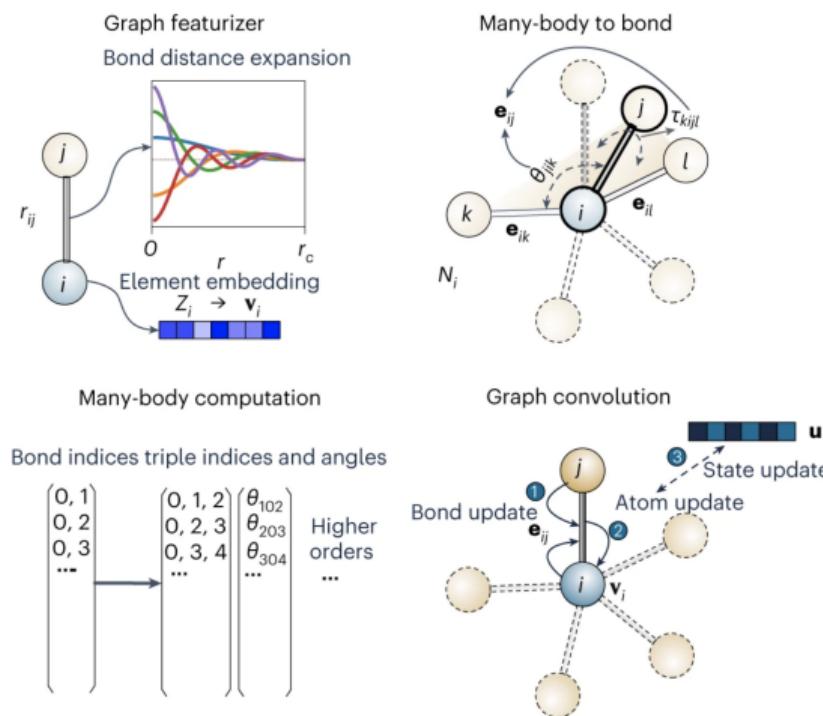


Figure: Schematic of the many-body graph potential [Chen and Ong, 2022]

O(3)-Equivariance & 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 accordingly.
- Need to combine node and edge features in a highly structured, prescriptive manner.

Clebsch-Gordon (CG) Tensor Product

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

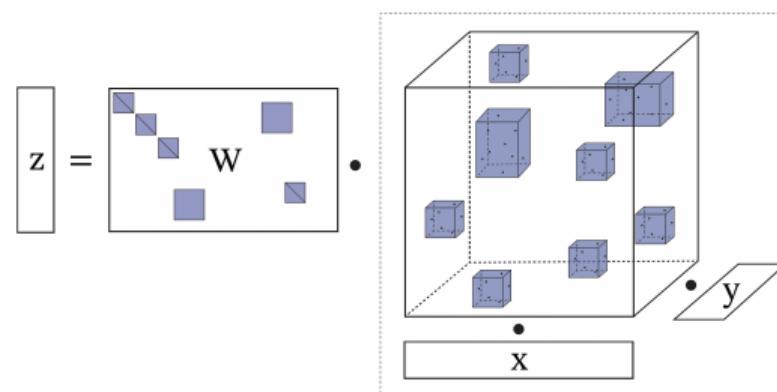


Figure: The CG tensor product, which contracts a block-sparse tensor \mathbf{P} with two dense vectors to output a new vector. Each blue block is itself sparse. [Bharadwaj et al., 2025]

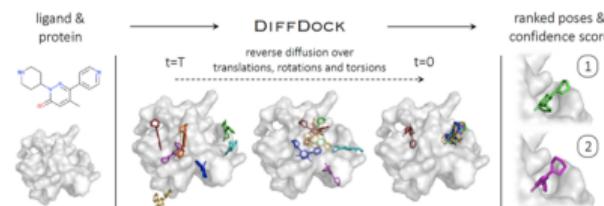
Bottlenecks to Scaling Equivariant Neural Networks

- ➊ Computational cost: the **CG tensor product** and **its derivatives** must be **evaluated millions of times** on large atomic datasets.
- ➋ **Memory bottlenecks** from computing large CG tensor products.
- ➌ Existing libraries launch **separate kernels** for each computation.

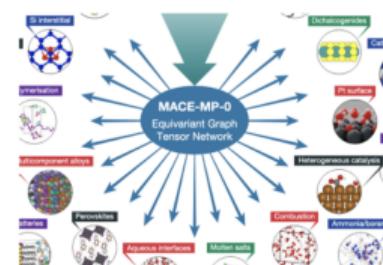
Existing Implementations



(a) Nequip [Bat+22]



(b) DiffDock [Cor+23]



(c) MACE-MP0 [Bat+24]

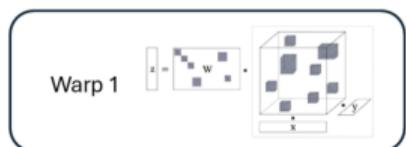
Solution: *OpenEquivariance*

OpenEquivariance is an open source **CUDA kernel generator** for the CG tensor product and its gradients. It achieves 10x speedup over the e3nn package, and 2x over NVIDIA cuEquivariance.

- Exploits sparse structure in the multilinear map tensor via JIT (Just-In-Time Compilation), register caching, and loop unrolling.
- Sparsity-aware forward and backward passes.
- Fuses CG tensor product with graph convolution.

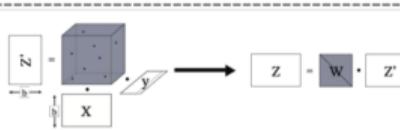
A Roadmap to Efficient CG Kernels

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, 0f3F080677;
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;
```



Warp-Level Forward Parallelism

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

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

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

Gradient Calculations

- Given $\frac{\partial E}{\partial z}$, want to compute $\frac{\partial E}{\partial x}$, $\frac{\partial E}{\partial y}$, $\frac{\partial E}{\partial W}$.
- Recompute forward pass to avoid memory traffic.
- Weight matrix handled through warp matrix multiplication.

Algorithm Subkernel C Warp-Level Backward

Require: $X \in \mathbb{R}^{b' \times (2\ell_x + 1)}$, $y \in \mathbb{R}^{2\ell_y + 1}$, $W \in \mathbb{R}^{b \times b'}$
Require: $G_Z \in \mathbb{R}^{b \times (2\ell_z + 1)}$, sparse tensor $\mathcal{P}^{(\ell_x, \ell_y, \ell_z)}$
Threads collaboratively compute $G'_Z = W^\top \cdot G_Z$

for $t = 1 \dots b'$ **do** ▷ Parallel over threads
 Load $x_{\text{reg}} = X[t, :]$, $y_{\text{reg}} = y$, $g'_{z\text{reg}} = G'_Z[t, :]$
 Initialize g_x reg, g_y reg, g_w reg, z reg to 0.

for $(i, j, k, v) \in \text{nz}(\mathcal{P}^{(\ell_x, \ell_y, \ell_z)})$ **do** ▷ Unroll via JIT
 g_x reg $[i] += v \cdot y_{\text{reg}}[j] \cdot g'_{z\text{reg}}[k]$
 g_y reg $[j] += v \cdot x_{\text{reg}}[i] \cdot g'_{z\text{reg}}[k]$
 $z_{\text{reg}}[k] += v \cdot x_{\text{reg}}[i] \cdot y_{\text{reg}}[j]$

Store g_y = warp-reduce(g_y reg)
Store $G_x[t, :] = g_x$ reg and $Z'[t, :] = z_{\text{reg}}$
Threads collaboratively compute $G_W = G_Z \cdot (Z')^\top$

Kernel Fusion for Graph Neural Networks (GNNs)

- Graph convolution has an SpMM memory access pattern.
- Fuse both kernels to save compute *and* memory.

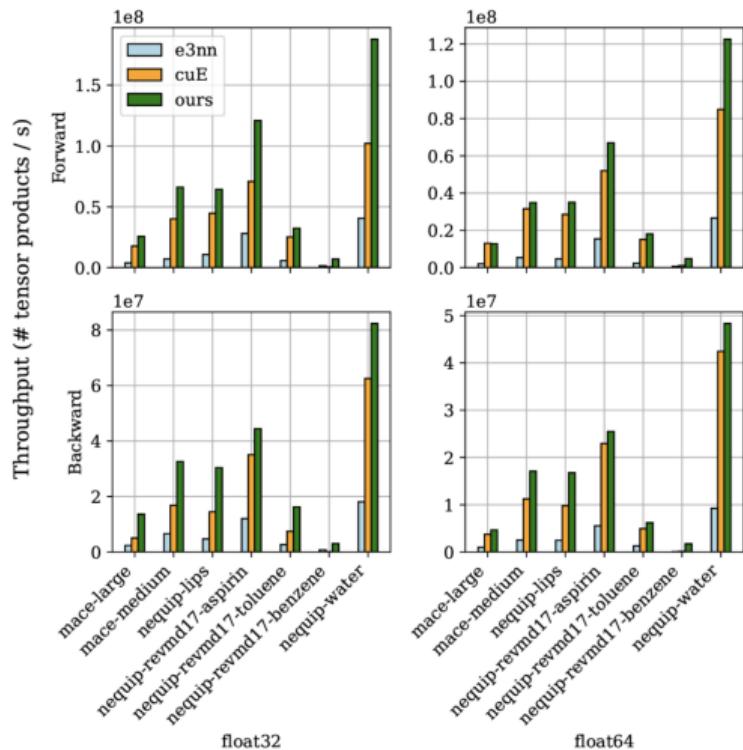
Algorithm Deterministic TP + Graph Convolution

```
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 segmenti ∈ schedule do
     $(s, t) = E [k] [0], E [k] [1]$ 
    Set  $z_{acc} = 0$                                 ▷ Parallel over Warps
    for  $b = 1 \dots |E|$  do
        Execute segment subkernel sequence
         $z_{acc} += z_{smem}$ 
        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.
             $z_{acc} = 0$ 
        Execute fixup kernel.
```

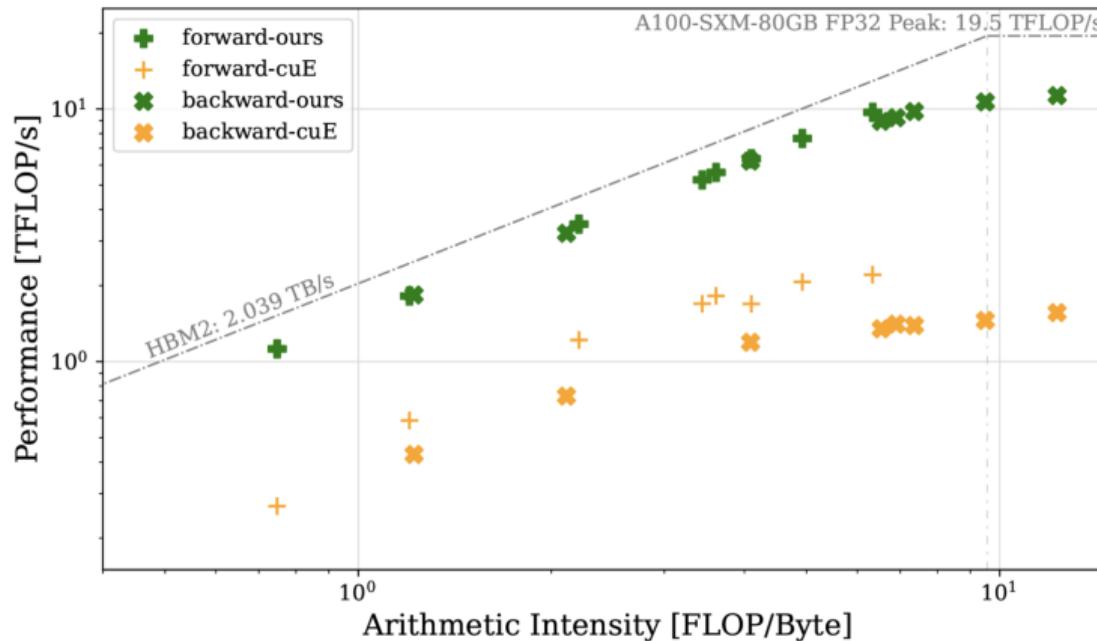
Experimental Results & Performance Gains

- ① Throughput without Kernel Fusion.
- ② Roofline Analysis.
- ③ Speedup with Kernel Fusion.
- ④ Acceleration of MACE Foundation Model.

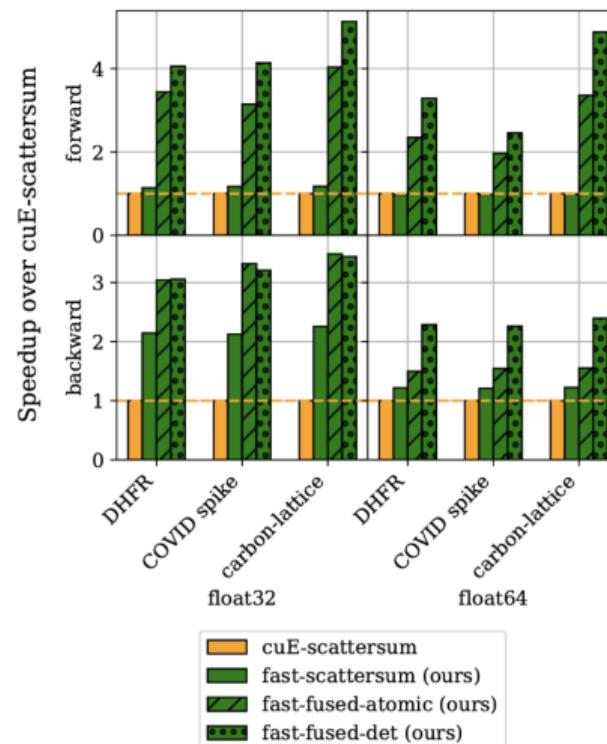
Throughput without Kernel Fusion



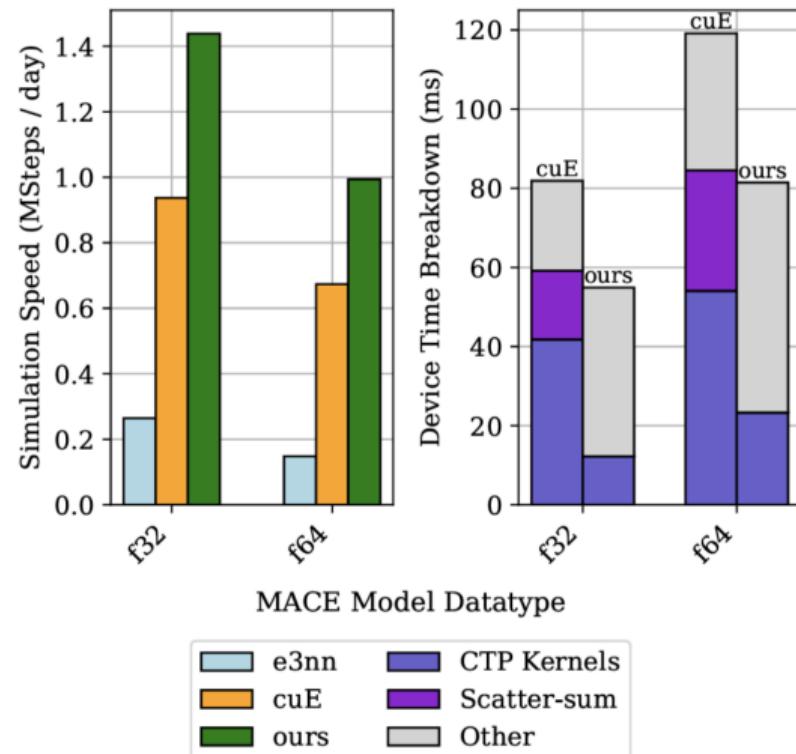
Roofline Analysis (F32)



Speedup with Kernel Fusion



Acceleration of MACE Foundation Model



Key Takeaways

- Sparse kernels achieve **consistently high performance** on the most common primitives used in equivariant deep neural networks.
- A **sparsity-aware, JIT-compiled** kernel is crucial.
- Future directions:
 - Reliance on the single-instruction multiple thread (SIMT) cores for FP32 and FP64 floating point arithmetic.
 - Expansion to specialized hardware for lower-precision computation.
 - Compile kernels to HIP for AMD GPUs.
 - Offer bindings for JAX and the Julia scientific computing language.

References

-  Bharadwaj, V., Glover, A., Buluc, A., and Demmel, J. (2025).
An efficient sparse kernel generator for $O(3)$ -equivariant deep networks.
-  Chen, C. and Ong, S. (2022).
A universal graph deep learning interatomic potential for the periodic table.
Nature Computational Science, 2:718–728.