



May 7th, 2025

Performance Portable Fast Ewald Summation

Kokkos User Group | Chicago, IL

Gabriel Kosmacher, Ziyu Du, Joar Bagge, George Biros



Predictive Engineering & Computational Science





 $\bullet\,$ Consider periodic N-body sums of the form

$$\boldsymbol{u}(\boldsymbol{x}_i) = \sum_{j=1}^{N_s} \sum_{\boldsymbol{p} \in \mathcal{P}} \boldsymbol{G}(\boldsymbol{x}_i - \boldsymbol{y}_j + \boldsymbol{p}) \boldsymbol{f}(\boldsymbol{y}_j) \quad i = 1, \dots, N_t.$$

• Ewald's method splits

$$oldsymbol{u}(oldsymbol{x}_i) = oldsymbol{u}^{\mathrm{N}}(oldsymbol{x}_i) + oldsymbol{u}^{\mathrm{F}}(oldsymbol{x}_i)$$

by setting $oldsymbol{G}(oldsymbol{r}) = oldsymbol{G}^{\mathrm{N}}(oldsymbol{r};\xi) + oldsymbol{G}^{\mathrm{F}}(oldsymbol{r};\xi).$

• For the electrostatic potential,

$$oldsymbol{G}(oldsymbol{r}) = rac{1}{\|oldsymbol{r}\|} = rac{ ext{erfc}(\xi\|oldsymbol{r}\|)}{\|oldsymbol{r}\|} + rac{ ext{erf}(\xi\|oldsymbol{r}\|)}{\|oldsymbol{r}\|} + rac{ ext{erf}(\xi\|oldsymbol{r}\|)}{\|oldsymbol{r}\|}.$$

 $\begin{array}{c} 10 \\ & & & \\ & &$

G. Kosmacher, ZD, JB, GB



 $\bullet\,$ Consider periodic N-body sums of the form

$$\boldsymbol{u}(\boldsymbol{x}_i) = \sum_{j=1}^{N_s} \sum_{\boldsymbol{p} \in \mathcal{P}} \boldsymbol{G}(\boldsymbol{x}_i - \boldsymbol{y}_j + \boldsymbol{p}) \boldsymbol{f}(\boldsymbol{y}_j) \quad i = 1, \dots, N_t.$$

• Ewald's method splits

$$\boldsymbol{u}(\boldsymbol{x}_i) = \boldsymbol{u}^{\mathrm{N}}(\boldsymbol{x}_i) + \boldsymbol{u}^{\mathrm{F}}(\boldsymbol{x}_i)$$

- by setting ${m G}({m r}) = {m G}^{
 m N}({m r};\xi) + {m G}^{
 m F}({m r};\xi).$
- For the electrostatic potential,

$$egin{aligned} egin{aligned} egi$$



G. Kosmacher, ZD, JB, GB



 $\bullet\,$ Consider periodic N-body sums of the form

$$\boldsymbol{u}(\boldsymbol{x}_i) = \sum_{j=1}^{N_s} \sum_{\boldsymbol{p} \in \mathcal{P}} \boldsymbol{G}(\boldsymbol{x}_i - \boldsymbol{y}_j + \boldsymbol{p}) \boldsymbol{f}(\boldsymbol{y}_j) \quad i = 1, \dots, N_t.$$

• Ewald's method splits

$$\boldsymbol{u}(\boldsymbol{x}_i) = \boldsymbol{u}^{\mathrm{N}}(\boldsymbol{x}_i) + \boldsymbol{u}^{\mathrm{F}}(\boldsymbol{x}_i)$$

by setting $oldsymbol{G}(oldsymbol{r}) = oldsymbol{G}^{\mathrm{N}}(oldsymbol{r};\xi) + oldsymbol{G}^{\mathrm{F}}(oldsymbol{r};\xi).$

• For the electrostatic potential,

$$egin{aligned} egin{aligned} egin{aligned} egin{aligned} egin{aligned} egin{aligned} egin{aligned} & egin{aligned} & \operatorname{erf}(\xi \|m{r}\|) \ & egin{aligned} & egin{aligned} & egin{aligned} & \operatorname{erf}(\xi \|m{r}\|) \ & egin{aligned} & egin{aligned} & egin{aligned} & egin{aligned} & \operatorname{erf}(\xi \|m{r}\|) \ & egin{aligned} &$$

G. Kosmacher, ZD, JB, GB





- Padded cell lists C^N accelerate **P2P**, **P2G**, and **G2P** by assigning each Kokkos team T to a cell β so each thread $t_x \in T$ processes spatially close data.
- + $\mathcal{C}^{\rm N}$ is constructed in two passes though an unsorted list of particles



GPU acceleration: Particle-to-Particle (P2P)



• For the stokes potential,

$$oldsymbol{G}^{\mathrm{N}}(oldsymbol{r}) = \left(oldsymbol{I} + rac{oldsymbol{r}\otimesoldsymbol{r}}{\|oldsymbol{r}\|^2}
ight) \left(\mathrm{erfc}(\xi\|oldsymbol{r}\|) + rac{\|oldsymbol{r}\| 2\xi\mathrm{e}^{-\xi^2\|oldsymbol{r}\|^2}}{\sqrt{\pi}}
ight) - oldsymbol{I}rac{4\xi\mathrm{e}^{-\xi^2\|oldsymbol{r}\|^2}}{\sqrt{\pi}}.$$

Table: FP64 functions on the NVIDIA H200 GPU, measured with NVIDIA Nsight Compute.

Operation	Millicycles	Flop-equivalents	Actual flops
DFMA	0.121	2	2
DADD	0.127	2.10	1
DMUL	0.127	2.10	1
DIV	2.11	35.0	15 (2)
RSQRT	1.82	30.2	8
EXP	3.07	50.8	30
ERFC	16.1	267	112

GPU acceleration: Particle-to-Particle (P2P)



• For the stokes potential,

$$oldsymbol{G}^{\mathrm{N}}(oldsymbol{r}) = \left(oldsymbol{I} + rac{oldsymbol{r}\otimesoldsymbol{r}}{\|oldsymbol{r}\|^2}
ight) \left(\mathrm{erfc}(\xi\|oldsymbol{r}\|) + rac{\|oldsymbol{r}\|^2 2\xi\mathrm{e}^{-\xi^2\|oldsymbol{r}\|^2}}{\sqrt{\pi}}
ight) - oldsymbol{I}rac{4\xi\mathrm{e}^{-\xi^2\|oldsymbol{r}\|^2}}{\sqrt{\pi}}.$$





The window $w(\mathbf{r}) = \sum_{i=1}^{d} w_0(r_i)$ is identical in each dimension. How can we use this structure to build an efficient parallel scheme?

P2G-BASE: Too many atomics! P2G-GRID: Redundant work! P2G-HYBRID: Best of both worlds?



The window $w(\mathbf{r}) = \sum_{i=1}^{d} w_0(r_i)$ is identical in each dimension. How can we use this structure to build an efficient parallel scheme?

P2G-BASE: Too many atomics!

P2G-GRID: Redundant work!

P2G-HYBRID: Best of both worlds?





The window $w(r) = \sum_{i=1}^{d} w_0(r_i)$ is identical in each dimension. How can we use this structure to build an efficient parallel scheme?

P2G-BASE: Too many atomics! P2G-GRID: Redundant work!





The window $w(r) = \sum_{i=1}^{d} w_0(r_i)$ is identical in each dimension. How can we use this structure to build an efficient parallel scheme?

P2G-BASE: Too many atomics! P2G-GRID: Redundant work! P2G-HYBRID: Best of both worlds?



Results: Performance Model





G. Kosmacher, ZD, JB, GB

May 7th, 2025 | Chicago, IL | #HPSFcon 5 / 8



Table: P2P results on an NVIDIA H200 GPU. Higher $\frac{N_s}{\text{us}}$ is better.

N	un at la a d	s = 256				
$I \mathbf{v}_t$	method	(b_x, b_y)	$\frac{N_t}{\mu s}$	Flops		
10^{6}	GM-1D	(64, 1)	14.8	71%		
	GM-2D	(1, 32)	11.8	56%		
	SM-1D	(32, 1)	15.2	72%		
	SM-2D	(64, 2)	12.4	59%		
$4 \cdot 10^{6}$	GM-1D	(64, 1)	15.3	73%		
	GM-2D	(8, 32)	11.9	57%		
	SM-1D	(32, 1)	15.5	74%		
	SM-2D	(64, 2)	12.4	59%		

Table: P2G results on an NVIDIA H200 GPU. Higher $\frac{N_s}{us}$ is better.

λŢ	method	$P = 8, E = 10^{-9}$			
$I \mathbb{V}_S$			Spd	Mops	
10^{6}	BASE	10.1		22%	
	GRID	11.9	$1.2 \times$		
	HYBRID	63.8	$6.3 \times$		
$4 \cdot 10^6$	BASE	8.5		18%	
	GRID	11.1	$1.3 \times$	24%	
	HYBRID	59.9	$7.1 \times$	28%	



Table: P2P results on an NVIDIA H200 GPU. Higher $\frac{N_a}{\mu s}$ is better.

N_t	method	s = 256				
				Flops		
10^{6}	GM-1D	(64, 1)	14.8	71%		
	GM-2D	(1, 32)	11.8	56%		
	SM-1D	(32, 1)	15.2	72%		
	SM-2D	(64, 2)	12.4	59%		
$4 \cdot 10^{6}$	GM-1D	(64, 1)	15.3	73%		
	GM-2D	(8, 32)	11.9	57%		
	SM-1D	(32, 1)	15.5	74%		
	SM-2D	(64, 2)	12.4	59%		

Table: P2G results on an NVIDIA H200 GPU. Higher $\frac{N_s}{us}$ is better.

λ	method	$P = 8, E = 10^{-9}$			
118		$\frac{N_s}{\mu s}$	Spd	Mops	
10^{6}	BASE	10.1	—	22%	
	GRID	11.9	$1.2 \times$	26%	
	HYBRID	63.8	$6.3 \times$	30%	
$4 \cdot 10^6$	BASE	8.5		18%	
	GRID	11.1	$1.3 \times$	24%	
	HYBRID	59.9	$7.1 \times$	28%	





Figure: Particle-to-particle performance portability (higher is better), s = 256.

G. Kosmacher, ZD, JB, GB





Figure: Particle-to-grid performance portability (higher is better). $E = 10^{-9}$, s = 224.

G. Kosmacher, ZD, JB, GB



- Introduced algorithms for spectral Stokes Ewald sum with performance analysis.
- Library implemented in PyKokkos, MPI results show preliminary overheads.
- Limitations (i.e., future work):
 - No single-precision tests.
 - CPU kernels and MPI communication not optimized.
 - Nonuniformity effects not studied.
 - Fix shared memory slowdowns on AMD GPUs.
 - No multi-GPU NUMA scaling.

Thank You!

Gabriel Kosmacher

gkosmacher@utexas.edu



- Introduced algorithms for spectral Stokes Ewald sum with performance analysis.
- Library implemented in PyKokkos, MPI results show preliminary overheads.
- Limitations (i.e., future work):
 - No single-precision tests.
 - CPU kernels and MPI communication not optimized.
 - Nonuniformity effects not studied.
 - Fix shared memory slowdowns on AMD GPUs.
 - No multi-GPU NUMA scaling.

Thank You!

Gabriel Kosmacher

gkosmacher@utexas.edu



• Particle-to-particle (P2P):

$$\boldsymbol{u}^{N}(\boldsymbol{x}_{i}) := \sum_{j \in [N_{s}]} \sum_{\boldsymbol{p} \in \mathcal{P}} \boldsymbol{G}^{N}(\boldsymbol{x}_{i} - \boldsymbol{y}_{j} + \boldsymbol{p}) \boldsymbol{f}(\boldsymbol{y}_{j}) \qquad \qquad \mathcal{O}(27sN_{t})$$

• Particle-to-grid (P2G):

$$\phi_h(g_\ell) := \sum_{j \in [N_s]} \sum_{p \in \mathcal{P}} w(g_\ell - y_j + p) f(y_j) \qquad \qquad \mathcal{O}(N_s P^3)$$

• Fourier Grid Convolution (FGC):

$$oldsymbol{v}_h(oldsymbol{g}_\ell) := extsf{IFFT}ig\{\widehat{oldsymbol{G}}^{ extsf{F}}(oldsymbol{\kappa}_\ell)]^{-2} extsf{FFT}ig\{oldsymbol{\phi}_h(oldsymbol{g}_\ell)ig\}ig\} \qquad \mathcal{O}(N_g\log N_g)$$

• Grid-to-particle (G2P):

$$\boldsymbol{u}_h^{\mathrm{F}}(\boldsymbol{x}_i) := h^3 \sum_{\ell \in [N_g]} \sum_{\boldsymbol{p} \in \mathcal{P}} w(\boldsymbol{x}_i - \boldsymbol{g}_\ell + \boldsymbol{p}) \boldsymbol{v}_h(\boldsymbol{g}_\ell) \qquad \qquad \mathcal{O}(N_t P^3)$$



0: parfor $x_i \in \beta$ do // Team gets β , TeamThreadRange over \boldsymbol{x}_i parfor $c \leftarrow 1, \ldots, \lceil s / |\alpha_c| \rceil$ do // VectorThreadRange over y_i 1: $\forall \boldsymbol{y}_i \in \alpha_c, \text{LOADSCRATCH}(\boldsymbol{y}_i, \boldsymbol{f}(\boldsymbol{y}_i))$ // GM/SM 2: TEAMBARRIER() 3: for $y_j \in \alpha_c$ do 4: $\mathbf{u}_{c} \leftarrow \mathbf{u}_{c} + \mathbf{G}^{N}(\mathbf{x}_{i} - \mathbf{u}_{i} + \mathbf{p}) \mathbf{f}(\mathbf{u}_{i})$ $// \mathbf{G}^{\mathrm{N}}$ is flop intensive! 5: 6: end for TEAMBARRIER() 7: 8: end parfor $\boldsymbol{u}_{h}^{N}(\boldsymbol{x}_{i}) \leftarrow \text{VECTORREDUCE}(\boldsymbol{u}_{c})$ Kokkos provided reduction 9: 10: end parfor

G. Kosmacher, ZD, JB, GB

Fast Ewald

May 7th, 2025 | Chicago, IL | #HPSFcon 2 / 5



0: parfor $y_i \in \beta$ do // Team gets β , TeamThreadRange over y_i Compute $w_1[l], w_2[l], w_3[l]$ for $l \leftarrow 0, \ldots, P-1$ and store in shared memory. 1: 2: end parfor 3: TEAMBARRIER() 4: parfor $t_x \leftarrow 0, \dots, (P/2)^3 - 1$ do // TeamThreadRange over $(P/2)^3$ threads for $\alpha \in \text{Colleagues}(\beta)$ do 5 $q_{\ell} \leftarrow \text{INDEX2GRID}(t_x, \alpha)$ // Assign q_{ℓ} to a thread 6. for $y_i \in \beta$ do 7: LOADSCRATCH $(w_1[(q_{\ell} - a_i)_1], w_2[(q_{\ell} - a_i)_2], w_3[(q_{\ell} - a_i)_3], f(y_i))$ 8: $\phi_{\text{loc}} \leftarrow \phi_{\text{loc}} + w_1 w_2 w_3 f(\mathbf{y}_i)$ 9: end for $10 \cdot$ ATOMICADD($\phi(q_{\ell}), \phi_{loc}$) $// \mathcal{O}(27N_s)$ atomics 11: 12: end for 13: end parfor

G. Kosmacher, ZD, JB, GB





Figure: Spectral Ewald run on multiple machines (lower is better). $E = 10^{-9}$, s = 224.

G. Kosmacher, ZD, JB, GB



Table: Runtime in *milliseconds* for multi-GPU weak scaling test with $N_{\rm GPU}$ NVIDIA H200 GPUs, MPI communication, and NVIDIA's cuFFTMp library. $E=10^{-9}$, s=224, $N=4\times10^6N_{\rm GPU}$.

Step	$N_{\rm GPU} =$	1 2	4	8	16
P2P	217.7	223.8	225.4	226.9	226.3
P2G	72.2	75.7	75.7	75.7	75.5
FFT	46.1	581.8	737.4	772.2	686.1
IFFT	1.5	24.0	28.6	30.7	26.1
G2P	39.9	40.6	40.6	40.4	40.7
MPI-SORT	5.6	22.7	23.2	36.4	61.7
MPI-GHOST-SOURCE	0.4	6.8	7.0	6.6	6.4
MPI-GHOST-GRID	0.2	3.2	5.2	4.7	4.6