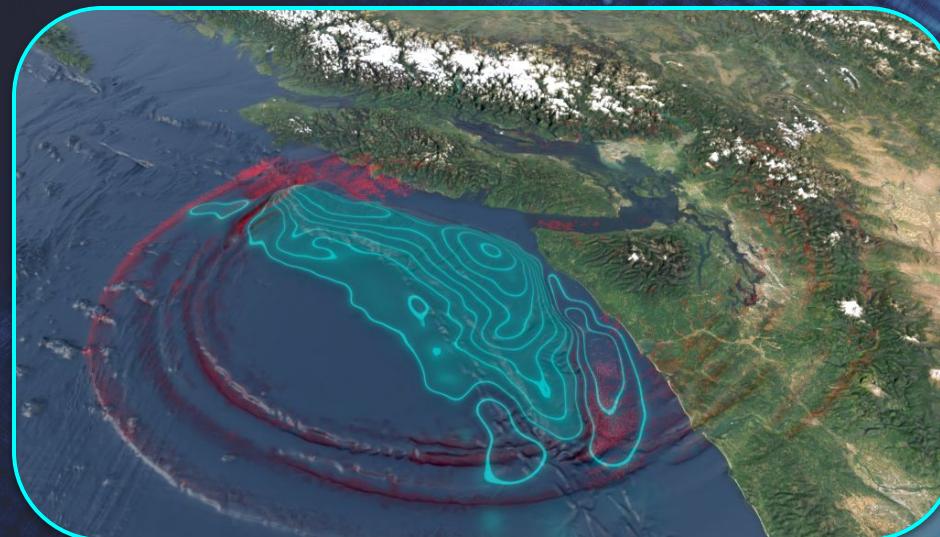




MFEM
Community
Workshop

A Guided Tour of MFEM GPU Kernel Optimization Techniques

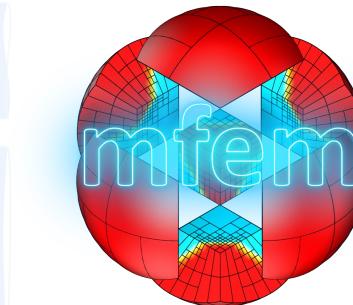
John Camier - LLNL Collaborator,
Veselin Dobrev, Tzanio Kolev - LLNL
Stefan Henneking - Oden Institute, The University of Texas at Austin
Jiqun Tu - NVIDIA



MFEM Workshop - GPU Kernel Optimizations Guided Tour

- Welcome students, new users & developers
- Exploring GPU kernel optimization strategies in MFEM

- Arbitrary order curvilinear **mesh** elements
- Arbitrary order **H1**, H(curl), H(div) and **L2** elements
- **Bilinear**/linear forms for: Galerkin, DG, etc.
- MPI-scalable assembly and linear solvers
- **GPU** acceleration on **AMD**, **NVIDIA** hardware
- Non-linear operators and non-linear solvers
- Explicit and implicit high-order time integration
- Integration with: hypre, SUNDIALS, SuperLU, PETSc, etc.



mfem.org v4.8 - Apr 2025



⇒ Real-time Bayesian inference at extreme scale:
A digital twin for tsunami early warning applied to the Cascadia subduction zone^[1]



Lawrence Livermore National Laboratory

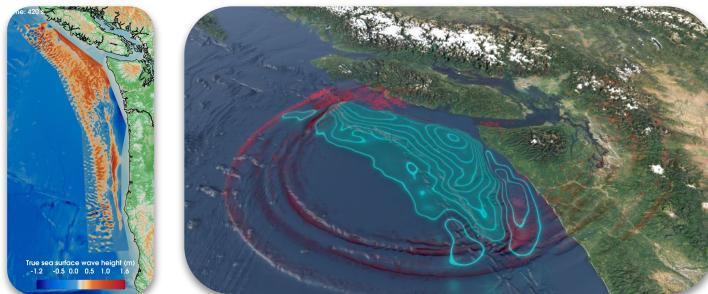
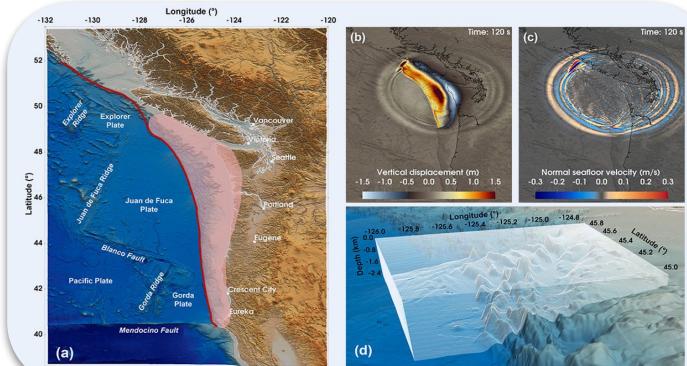
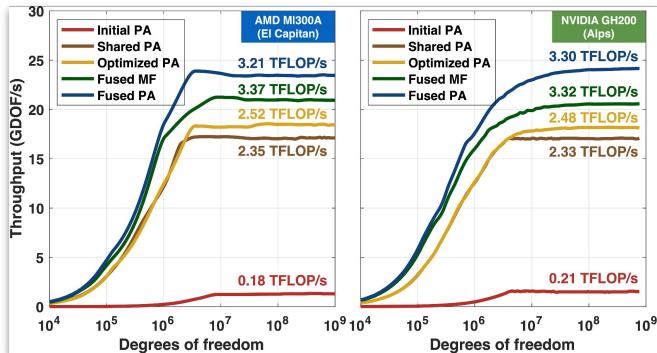


Application - A Digital Twin For Tsunami Early Warning

- Important and challenging problem
- Forecast wave heights or onshore inundation
- Produce better early warning systems for tsunamis

$$\begin{pmatrix} A \begin{bmatrix} \vec{u} \\ p \end{bmatrix}, \begin{bmatrix} \vec{\tau} \\ v \end{bmatrix} \end{pmatrix} := \begin{bmatrix} 0 & (\nabla p, \vec{\tau}) \\ -(\vec{u}, \nabla v) & \langle Z^{-1}p, v \rangle_{\partial\Omega_a} \end{bmatrix}$$

- Problem size \Rightarrow memory optimizations
- Key kernels \Rightarrow performance optimizations

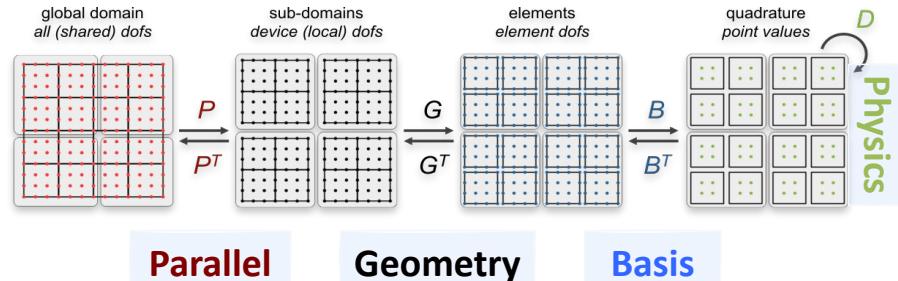
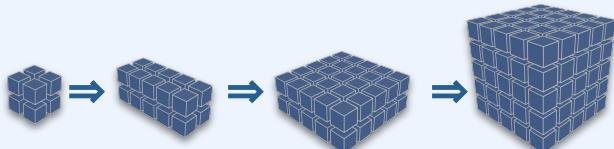


MFEM Operator Decomposition for GPU Kernels

- Partial Assembled HO Finite Element Operators
 - $A = P^T G^T B^T D B G P$
- Optimal memory, near-optimal FLOPs
- Matrix free: no assembly of the full matrix A

- GPU Kernel optimization: focus on $G^T B^T D B G$

- B :

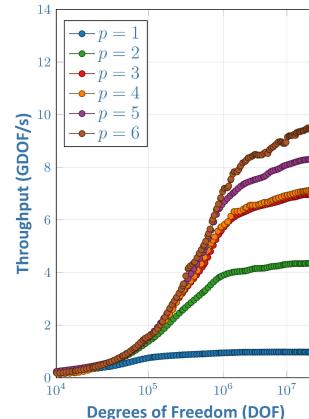


Parallel

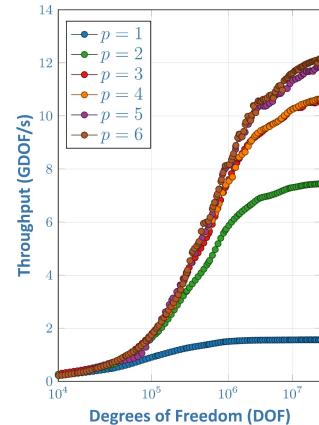
Geometry

Basis

MFEM BP1 (standard) @ GH200



MFEM BP1 (atomics) @ GH200



MFEM GPU Kernel Overview

Integration

MFEM's spatial integrations are performed in the usual finite element manner by first splitting the spatial domain into a collection of non-overlapping "elements" which cover the domain. This is usually referred to as the "mesh". An integral can then be computed separately in each element and the results added together:

$$\int_Q f(x) d\Omega = \sum_i \int_{\Omega_i} f(x) d\Omega$$

Where Ω is the full domain and Ω_i is the domain of the i -th element. In MFEM this sum over elements is performed in classes such as the `BilinearForm` or `LinearForm` and their parallel counterparts.

Bilinear Form Integrators

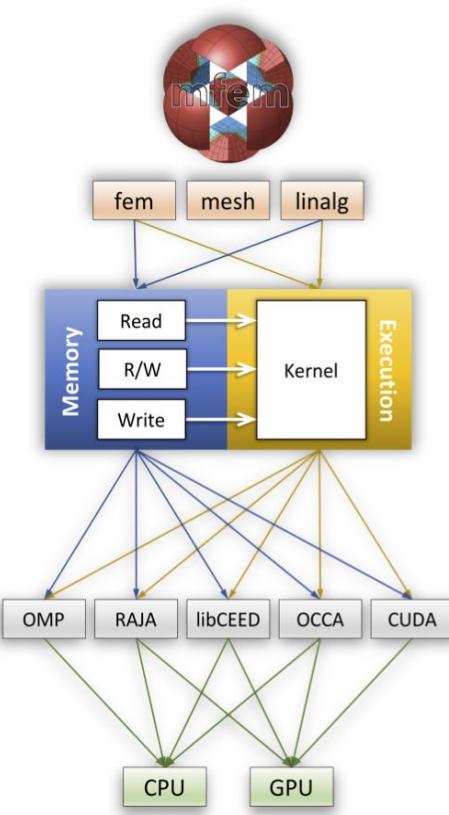
Bilinear form integrators are at the heart of any finite element method, they are used to compute the integrals of products of basis functions over individual mesh elements (or sometimes over edges or faces). Typically each element is contained in the support of several basis functions of both the domain and range spaces, therefore bilinear integrators simultaneously compute the integrals of all combinations of the relevant basis functions from the domain and range spaces. This produces a two dimensional array of results that are arranged into a small dense matrix of integral values called a local element (stiffness) matrix.

To put this another way, the `BilinearForm` class builds a global, sparse, finite element matrix, `glb_mat`, by performing the outer loop in the following pseudocode snippet whereas the `BilinearFormIntegrator` class performs the nested inner loops to compute the dense local element matrix, `loc_mat`.

```
for each elem in elements
    loc_mat = 0.0
    for each pt in quadrature_points
        for each u_j in elem
            for each v_i in elem
                loc_mat(i,j) += w(pt) * u_j(pt) v_i(pt)
        end
    end
    glb_mat += loc_mat
end
```

Mixed Operators

Class Name	Domain	Range	Coeff.	Dimension	Operator
VectorDivergenceIntegrator	H1 ^d , L2 ^d	H1, L2	S	1D, 2D, 3D	($\lambda \nabla \cdot \vec{u}, \vec{v}$)
GradientIntegrator	H1	H1 ^d , L2 ^d	S	1D, 2D, 3D	($\lambda \nabla u, \vec{v}$)



GPU support in MFEM

MFEM relies mainly on two features for running algorithms on devices such as GPUs:

- The memory manager handles transparently the moving of data between the host (CPU) and the device (e.g. GPU),
- The `mfem::forall` function to abstract `for` loops to parallelize the execution on an arbitrary device.

```
Vector u;
Vector v;
// ...
const auto u_data = u.Read(); // Express the intent to read u
auto v_data = v.ReadWrite(); // Express the intent to read and write v

// Abstract the loop: forall(u.Size(), [=] MFEM_HOST_DEVICE (int i)
{
    v_data[i] *= u_data[i]; // This block of code is executed on the chosen device
};
```

Supported Integrators

native MFEM
Mass Integrator
Vector Mass Integrator
Vector FE Mass Integrator
Convection Integrator
Non-linear Convection Integrator
Diffusion Integrator
Vector Diffusion Integrator
DGTrace Integrator
Mixed Vector Gradient Integrator
Mixed Vector Curl Integrator
Mixed Vector Weak Curl Integrator
Gradient Integrator
Vector Divergence Integrator
Vector C Div Divergence Integrator
Curl Curl Integrator
Div Div Integrator

- Memory: input/outputs data management
- Execution: outer/inner forall loops
- Kernel: Integrator, G^TB^TDBG



Initial PA - Integrators Implementation

- Mixed integrator: B_{test} , B_{trial}
 - Fixed order: $\{4_{\text{test}}, 5_{\text{trial}}\}$ dofs, 5 at quadrature points
 - CPU development, GPU portable
 - Extract: PA setup, Q function

Writing Custom Integrators

Element-wise integration arises in various places in the finite element method. A square and rectangular bilinear form operators, linear functionals, and the calcul

Type	Primary Function Needing Implementation
Square Operators	<code>BilinearFormIntegrator::AssembleElementMat</code>
Rectangular Operators	<code>BilinearFormIntegrator::AssembleElementMat</code>
Linear Functionals	<code>LinearFormIntegrator::AssembleRHSElementVec</code>

Name	C++ Expression	Formula
Jacobian Matrix	const DenseMatrix &J = Trans.Jacobian()	$J_{ij} = \frac{\partial f_i}{\partial x_j}$
Jacobian Determinant	double detJ = Trans.Weight()	$\det(J)$
Inverse Jacobian	const DenseMatrix &InvJ = Trans.InverseJacobian()	J^{-1}
Adjugate Jacobian	const DenseMatrix &adjJ = Trans.AdjugateJacobian()	$\det(J) J^{-1}$

```

void BuildIntegrator:AssembleElementMatrix()
{
    const FiniteElement &trial_fe, const FiniteElement &test_fe,
    ElementTransformation &trans, DenseMatrix &elmat{};

    dim = test_fe.GetDim();
    int trial_dof = trial_fe.GetDof(0); test_dof = test_fe.GetDof(0);
    real_t c;
    Vector d_col;

    dshape.SetSize(trial_dof, dim); gshape.SetSize(trial_dof, dim);
    Jadj.SetSize(dim), shape.SetSize(test_dof);
    elmat.SetSize(dim * test_dof, trial_dof);

    const IntegrationRule &ir = GetIntegrationRule(trial_fe, test_fe, Trans);
    weight = 8.0;
    elmat_comp.SetSize(test_dof, trial_dof);

    for (int i = 0; i < ir.GetDofPoints(); i++)
    {
        const IntegrationPoint &ip = ir->IntPoint(i);
        Trans.SelPoint(ip);

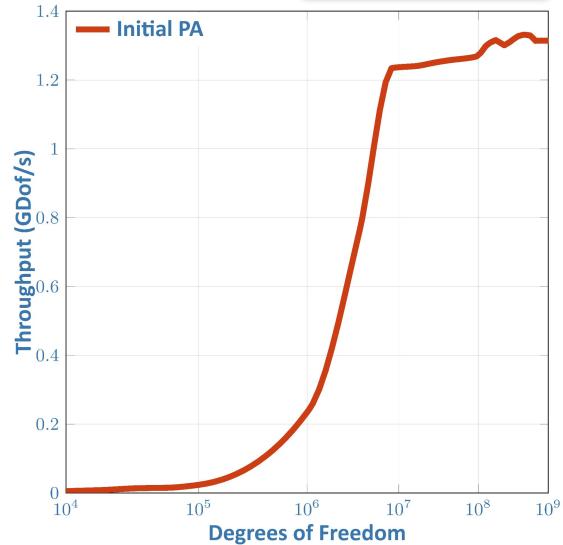
        CalcJugated(trans.Jacobian(), Jadj);
        test_fe.CalcPhiShape(trans, shape), trial_fe.CalcDShape(ip, dshape);
        MultiShape(Jadj, gshape);

        c = ip.weight;
        if (0) { c = 0; Eval(Trans, ip); }
        shape *= c;

        for (int d = 0; d < dim; ++d)
        {
            gshape.GetColumnForElement(d, d_col);
            MultiShape(gshape, d_col, elmat_comp);
            for (int jj = 0; jj < dim; ++jj)
                for (int ii = 0; ii < test_dof; ++ii)
                {
                    elmat[d * test_dof + ii, jj] += elmat_comp(ii, jj);
                }
        }
    }
}

```

AMD MI300A



GPU Kernel Optimizations - Profiling and Benchmarking

- Using Google benchmark: agile development, CPU & GPU timings
- tests/benchmarks examples

Benchmark	Time	CPU	Iterations	Dofs	MDof/s	p	version
WaveOp/1/4/160	333 ms	333 ms	2	16.4613M	49.4931/s	4	1
WaveOp/8/4/160	305 ms	305 ms	2	16.4613M	53.9399/s	4	8
WaveOp/3/4/160	187 ms	187 ms	4	16.4613M	88.1353/s	4	3
WaveOp/4/4/160	144 ms	144 ms	5	16.4613M	114.076/s	4	4
WaveOp/6/4/160	174 ms	174 ms	4	16.4613M	94.6562/s	4	6
WaveOp/9/4/160	473 ms	473 ms	2	16.4613M	34.799/s	4	9
WaveOp/13/4/160	242 ms	242 ms	3	16.4613M	67.9619/s	4	13



Apple M2 Pro

Benchmark	Time	CPU	Iterations	Dofs	MDof/s	p	version
WaveOp/1/4/160	17.4 ms	17.1 ms	41	16.4613M	960.653/s	4	1
WaveOp/8/4/160	1.83 ms	1.82 ms	385	16.4613M	9.06081K/s	4	8
WaveOp/3/4/160	0.897 ms	0.897 ms	773	16.4613M	18.3594K/s	4	3
WaveOp/4/4/160	0.703 ms	0.700 ms	983	16.4613M	23.5143K/s	4	4
WaveOp/6/4/160	1.20 ms	1.19 ms	586	16.4613M	13.8101K/s	4	6
WaveOp/9/4/160	1.07 ms	1.06 ms	656	16.4613M	15.531K/s	4	9
WaveOp/13/4/160	0.818 ms	0.810 ms	860	16.4613M	20.316K/s	4	13

AMD MI300A

- **-Rpass-analysis=kernel-resource-usage**
[S,V,A]GPRs, ScratchSize, Occupancy, LDS Size
- **--ptxas-options=-v**

void mfem::HipKernel1D<mfem::RK4... void mfem::hip::HipKernel3D<128, 1, mfem::internal::SmemPAGra... void mfem::HipKernel... void mf... void mfem::hip::HipKernel3D<128, 1, mfem::internal::SmemPA... v void mfem::HipKerne...

5M elements, 1.3B dofs

- rocprofv3 & <https://ui.perfetto.dev>
- NVIDIA Nsight Systems/Compute



Shared Memory PA Kernel Optimizations

void mfem::HipKernel1D<mfem::RK4...

void mfem::hip::HipKernel3D<128, 1, mfem::internal::SmemPAgra...

void mfem::HipKernel...

void mfem::hip::HipKernel3D<128, 1, mfem::internal::SmemPA...

y void mfem::HipKerne...

```

template <int D1DR, int D1DE, int QID>
void SmemPAgradientApplyTranspose3D<const int NE,
    const Array<real_t> &trial_bt,
    const Array<real_t> &trial_gt,
    const Array<real_t> &test_bt,
    const Vector<6d>
    const Vector<6x> Vector<6y>
    real_t &a, real_t &b>
{
    constexpr int M01 = D1DE + D1DR ? D1DE : D1DR;

    const auto D_r_t = Reshape(trial_bt.Read(0), D1DR, QID);
    const auto gr_r_t = Reshape(trial_gt.Read(0), D1DR, QID);
    const auto br_r_t = Reshape(test_bt.Read(0), D1DE, QID);

    const auto D_g = Reshape(d.Read(0), QID * QID * QID, 3, NE);
    const auto X_g = Reshape(x.Read(0), D1DE, D1DE, D1DE, NE, 3);
    const auto Y_g = Reshape(y.Read(0), D1DR, D1DR, D1DR, NE);
    auto Z_g = Reshape(z.Read(0), D1DR, D1DR, D1DR, NE);

    mfem::forall_3D<NE, QID, QID, QID, [-] MFEM_HOST_DEVICE(int e)
    {
        MFEM_SHARED real_t_BG[2](QID * M01);
        MFEM_SHARED real_t_sm[3](QID * QID * QID), sm1[3](QID * QID * QID);

        kernels::internal::LoadX0(M01, QID * e, D1DE, X_g, sm1);
        kernels::internal::LoadY0(M01, QID * e, D1DR, Y_g, sm1);
        kernels::internal::LoadZ0(M01, QID * e, D1DR, Z_g, sm1);

        kernels::internal::EvalX(M01, QID * QID, QID, BG[0], sm0, sm1);
        kernels::internal::EvalY(M01, QID * QID, QID, BG[0], sm0, sm1);
        kernels::internal::EvalZ(M01, QID * QID, QID, BG[0], sm0, sm1);

        MFEM_FOREACH_THREAD(qz, z, QID)
        {
            MFEM_FOREACH_THREAD(qy, y, QID)
            {
                MFEM_FOREACH_THREAD(qx, x, QID)
                {
                    real_t G13;
                    const int q = qy + qx * QID + QID * QID;
                    kernels::internal::PushEval0D(QID, QD, qx, qy, qz, sm1, G1);
                    A(0) = (D(q, 0, 0, e) * G13) + (D(q, 0, 1, e) * G11) + (D(q, 2, 0, e) * G23);
                    A(1) = (D(q, 0, 0, e) * G13) + (D(q, 1, 0, e) * G11) + (D(q, 1, 1, e) * G21);
                    A(2) = (D(q, 0, 0, e) * G13) + (D(q, 1, 0, e) * G11) + (D(q, 2, 0, e) * G22);
                    kernels::internal::PushGrad0D(QID, QD, qx, qy, qz, A, sm1);

                    kernels::internal::LoadBGt(M01, QID * D1DR, QD, br_r_t, gr_r_t, BG);
                    kernels::internal::Grad2t(M01, QID * D1DR, QD, BG, sm0, sm1);
                    kernels::internal::GradYt(M01, QID * D1DR, QD, BG, sm1, sm0);
                    kernels::internal::GradZt(M01, QID * D1DR, QD, BG, sm0, sm0);
                }
            }
        }
        MFEM_SYNC_THREAD;
    }
}

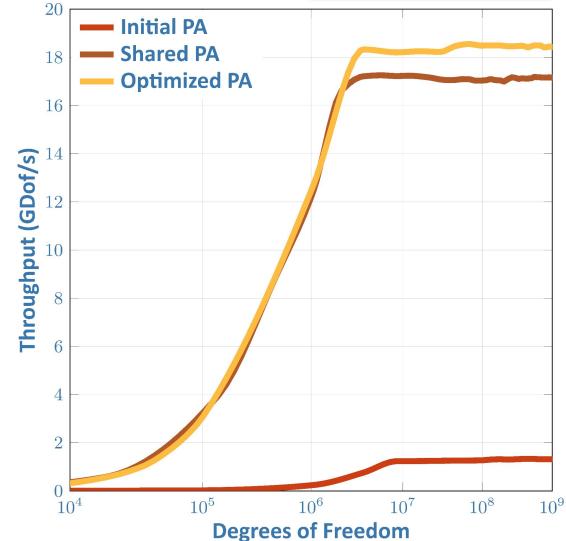
```

B^T_{trial} D^T B_{test}

G^T_{trial} G_{trial}

B^T_{test} D B^T_{trial}

AMD MI300A



- G^T, G : Local to Element (L to E) vectors handled by MFEM
- "Optimized PA" reached by fixing the launch bounds



Lawrence Livermore National Laboratory



NNSA
National Nuclear Security Administration

Application - HPC Context & Algorithms

```
void mfem::HipKernel1D<mfem::RK4...
```

```
void mfem::hip::HipKernel3D<128, 1, mfem::internal::SmemPAGra...
```

```
void mfem::HipKernel...
```

```
void mf...
```

```
void mfem::hip::HipKernel3D<128, 1, mfem::internal::SmemPAGra...
```

```
y void mfem::HipKernel...
```

B^T_{trial} D^T_{trial} B_{test}

G^T_{trial}

G_{trial}

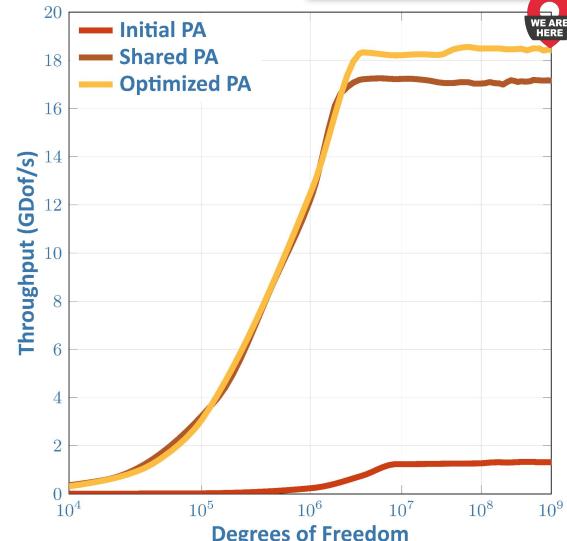
B^T_{test} D_{test} B^T_{trial}

AMD MI300A

- Vector sizes: $\mathcal{O}(\text{GB})$
- Multiple read of the same data should be avoided
- Overall optimization fusion pass

- G^T, G : replaced by indirection + atomics

$B^T_{\text{test}} D_{\text{test}} B^T_{\text{trial}} G_{\text{trial}} G^T_{\text{trial}} B^T_{\text{trial}} D^T B_{\text{test}} X_{u,p} \Rightarrow B^T_{\text{test}} B^T_{\text{trial}} D D^T B_{\text{trial}} B_{\text{test}} X_{u,p}$



```
void mfem::HipKernel1D<mfem::RK4Solv...
```

```
void mfem::hip::HipKernel3D<128, 1, mfem::internal::SmemPAGradientApplyTranspose3DRtRMult<5, 4, 5, 128, 1>(mfem::FiniteElementSpace const*, mfem...
```

```
y void mfem::HipKernel1D...
```

B^T_{trial} B^T_{test} $D D^T B_{\text{trial}}$ B_{test}

Fused PA Kernel Optimizations

- Reduced memory access: PA data read once

- Challenges:
 - register pressure
 - increased complexity
 - shared memory usage

- PA data uses **1/3** of the memory, w/:
 - avoiding caching large vectors
 - recomputing on-the-fly some values
 - reusing temporary vectors from RK4

$$\mathbf{B}_{\text{trial}}^T \mathbf{B}_{\text{test}}^T \mathbf{D} \mathbf{D}^T \mathbf{B}_{\text{trial}} \mathbf{B}_{\text{test}}$$

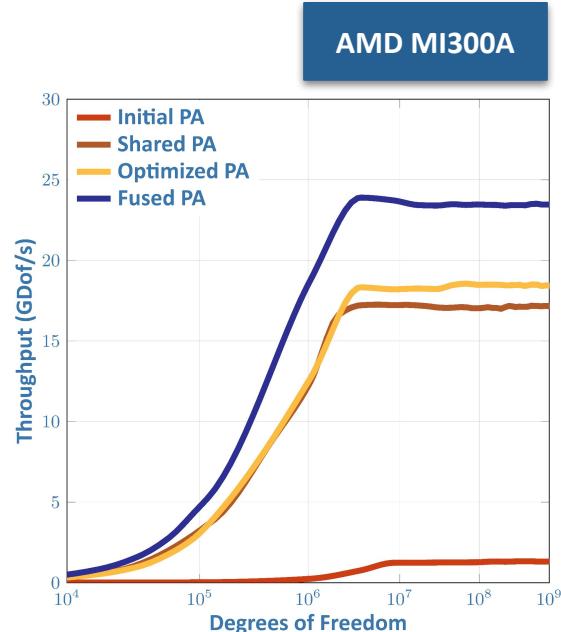
void mfem::HipKernel1D<mfem::RK4Solv...

```
void mfem::hip::HipKernel3D<128, 1, mfem::SmemPAGradientApplyTranspose3DRtRMult<5, 4, 5, 128, 1>(mfem::FiniteElementSpace const*, mfem...
```

void mfem::HipKernel1D<mfem::RK4...

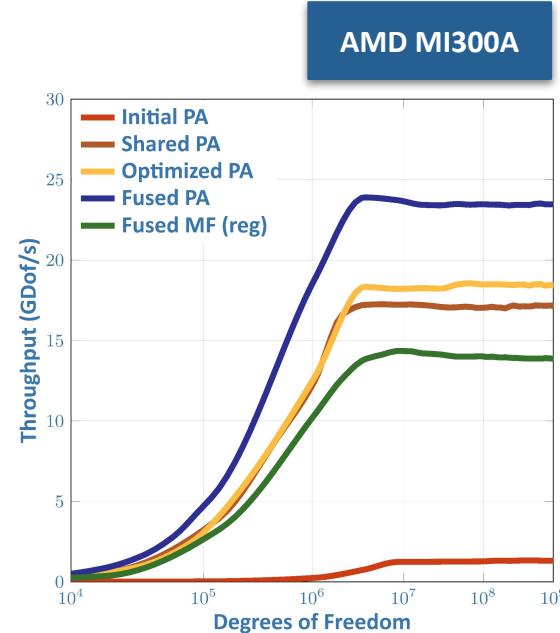
void mfem::HipKernel...

void mfem::hip::HipKernel3D<128, 1, mfem::internal::SmemPA...
mf...



Fused MF Kernel Optimizations - 1/2

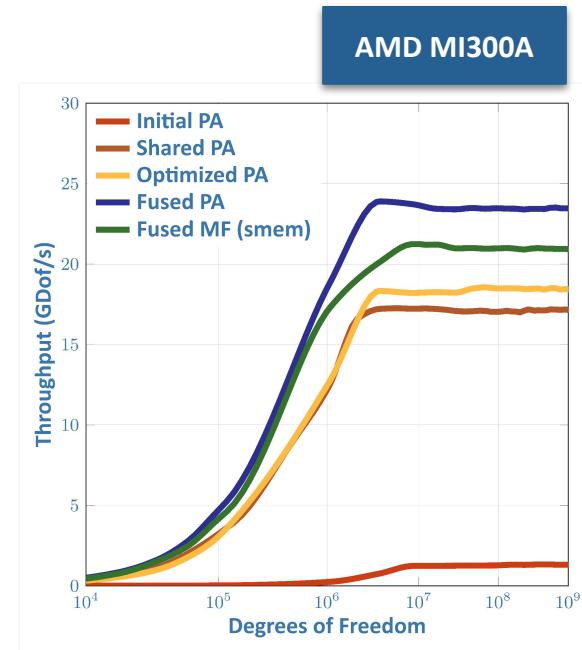
- No PA data stored at Quadrature points
 - Extra input vectors & computations
 - Indirections, basis arrays
 - Mesh coordinates: used for 'setup'
 - Sum factorisation 3D vector grad basis
 - Multiple implementations
 - `smem`: default, with 3D block of smem & thread
 - `reqs`: less shared mem. 2D thread blocks



```
void mfem::HipKernel1D<m... v void mfem::hip::HipKernel3D<128, 1, mfem::SmemMFGGradApplyT3DRtRMult<5, 4, 5, 2, 5, 128, 1>(mfem::FiniteElementSpace const*, mfem::ElementRestriction const*, int, ... void mfem::HipKernel1D<mfem::RK4Solv... vo... void mfem::hip::HipKernel3D<128, 1, mfem::SmemPAGradientApplyTranspose3DRtRMult<5, 4, 5, 128, 1>(mfem::FiniteElementSpace const*, mfem... v void mfem::HipKernel1D... void mfem::HipKernel1D<mfem::RK4... void mfem::hip::HipKernel3D<128, 1, mfem::internal::SmemPAGrad... void mfem::HipKernel1D<mfem::RK4... void mfem::hip::HipKernel3D<128, 1, mfem::internal::SmemPA... v void mfem::HipKernel1D<mfem::RK4...
```

Fused MF Kernel Optimizations - 2/2

- Increasing the occupancy: number of wavefronts
- Use compiler output:
 - 170 max VREG
 - 3 waves \Rightarrow 1638 maximum fp64
- Reducing register usage:
 - FORALL_DIRECT
- Reducing the amount of shared memory
 - move B_{trial} , B_{test} data to constant memory
 - shuffle/re-use vector grad computation



```
void mfem::HipKernel1D<mfem::RK4Solv... vo... void mfem::hip::HipKernel3D<128, 1, mfem::SmemPAGradientApplyTranspose3DRtRMult<5, 4, 5, 128, 1>(mfem::FiniteElementSpace const*, mfem::ElementRestriction const*, mfem::SmemPA*); v... void mfem::HipKernel1D...
void mfem::HipKernel1D<mfem::RK4Solv... v... void mfem::hip::HipKernel3D<128, 1, mfem::SmemMFGGradApplyT3DRtRMult_amd<5, 4, 5, 2, 5, 128, 1>(mfem::FiniteElementSpace const*, mfem::ElementRestriction const*, mfem::SmemPA*); v... void mfem::HipKernel1D...
void mfem::HipKernel1D<mfem::RK4... void mfem::hip::HipKernel3D<128, 1, mfem::internal::SmemPAGra... void mfem::HipKernel... void mf... void mfem::hip::HipKernel3D<128, 1, mfem::internal::SmemPA... v... void mfem::HipKerne...
```

Shared Memory, Fused PA & Fused MF NVIDIA Kernel Optimizations

```

Template<int DIM>
Int_DUNE::Int_QID, Int_MID = D3MH<
void SmechAdvectionType::TransposeAndScale(
    const FiniteElementSpace &trial_Fes, const ElementRestriction &trial_R,
    const Vector &coeffs)
{
    Arrayreal<D, dtrial_D, Arrayreal<D, dtrial_pt, Arrayreal<D, dttest_D,
    const Vector &dt, Vector &dt,
    const Vector &dtet, Vector &dtet,
    const real_t &realt, const real_t &dtet,
    const Vector &realt, Vector &dtet,
    Vector &realt_dt, Vector &dtet_dt,
    Vector &realt_dtet, Vector &dtet_dtet,
    Vector &realt_dtet_dt, Vector &dtet_dtet_dt);
    Arrayreal<D, dtrial_D, Arrayreal<D, dtrial_pt, Arrayreal<D, dttest_D,
    const Vector &dt, Vector &dt,
    const Vector &dtet, Vector &dtet,
    const real_t &realt, const real_t &dtet,
    const Vector &realt, Vector &dtet,
    Vector &realt_dt, Vector &dtet_dt,
    Vector &realt_dtet, Vector &dtet_dtet,
    Vector &realt_dtet_dt, Vector &dtet_dtet_dt);
}

```

SHARED

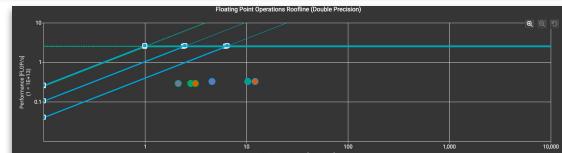
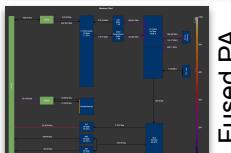
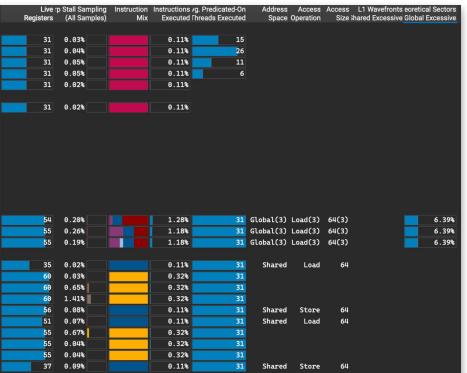
B_{test}
B_{trial}

D^T
D

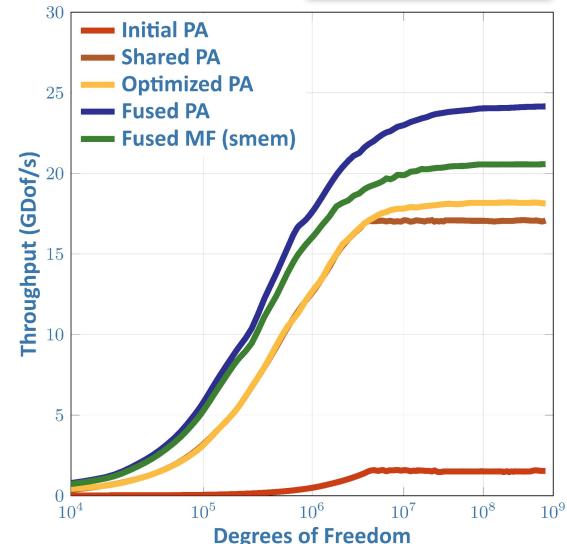
$$\begin{matrix} \mathbf{B}^T_{\text{trial}} \\ \mathbf{B}^T_{\text{test}} \end{matrix}$$

- Switch seamlessly to NVIDIA hardware
 - Resilient to the different *versions*
 - ✓ Shared memory bound kernel

5

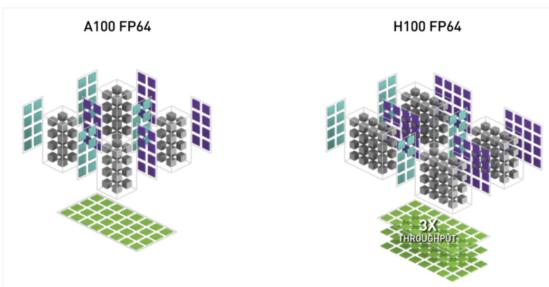


NVIDIA GH200

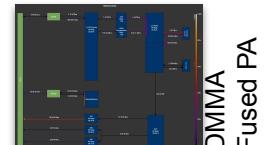
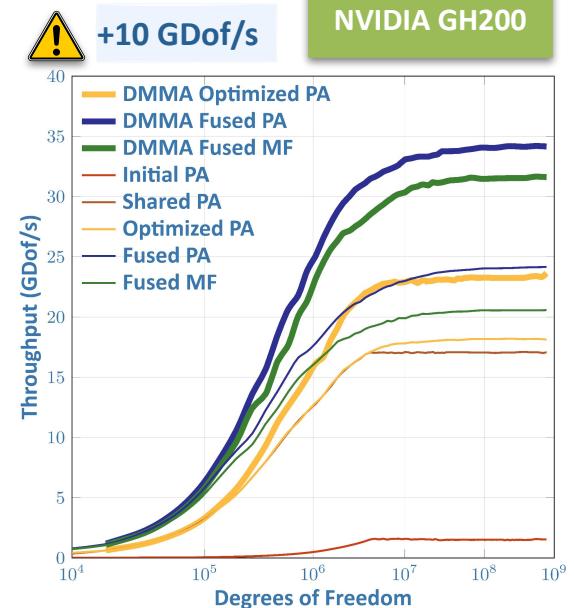


Unlocking Next-Level Performance Opportunities

- Jiqun Tu added Tensor Core based contractions
 - M-by-N-by-K warp-synchronous collectives

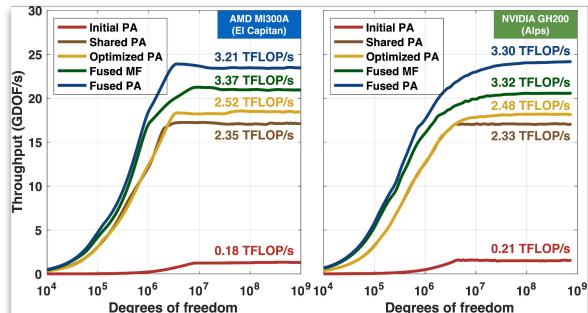
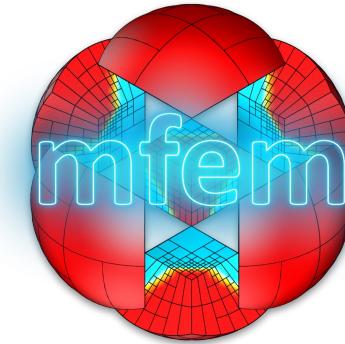


- 4th generation Matrix Multiply-Add (MMA)
 - $D = op(A, B) + C$
 - \Rightarrow Higher throughput
 - \Rightarrow More efficient way to share data
 - For shared memory bound kernels \Rightarrow speed



Conclusion

- Practical insights for enhancing FE HPC computations
 - Contributions are welcome!
-
- Holistic Kernel Fusion Approach
 - Not only limited to kernel launch overhead
 - Re-use data, avoid in-&-out data transfers
-
- WIP tensor contraction API to support:
 - Low vs. high order algorithms
 - Arbitrary number of arguments for ∂ FEM



mfem.org



[1] Henneking, Stefan, et al., Real-time Bayesian inference at extreme scale: A digital twin for tsunami early warning applied to the Cascadia subduction zone



Disclaimer

This document was prepared as an account of work sponsored by an agency of the United States government. Neither the United States government nor Lawrence Livermore National Security, LLC, nor any of their employees makes any warranty, expressed or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States government or Lawrence Livermore National Security, LLC. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States government or Lawrence Livermore National Security, LLC, and shall not be used for advertising or product endorsement purposes.