# libCEED GPU Strategy

Jeremy L Thompson

University of Colorado Boulder

*jeremy@jeremylt.org*

11 Feb 2025

# libCEED, Ratel, and HONEE Team



libCEED Repo: https://github.com/CEED/libCEED
HONEE Repo: https://gitlab.com/phypid/HONEE
Ratel Repo: https://gitlab.com/micromorph/ratel

Developers: **Zach Atkins**, Jed Brown, Fabio Di Gioacchino, Leila Ghaffari,
Kenneth Jansen, **Rezgar Shakeri**, James Wright,
**Jeremy L Thompson**

# Overview

## ECP Roots

- libCEED + PETSc projects follow from ECP CEED work

- libCEED provides high-performance operator evaluation

- libCEED provides CUDA, ROCm, and SYCL support

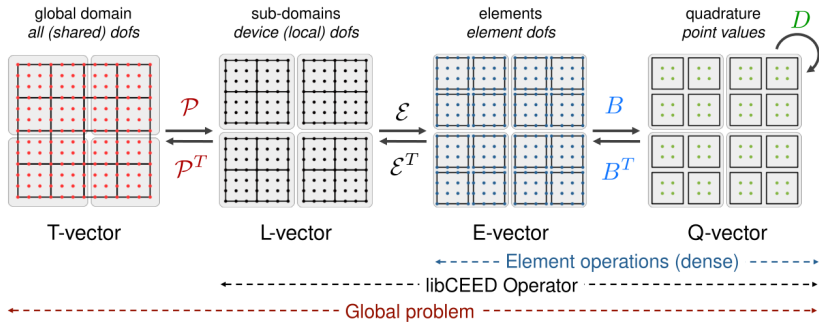- PETSc provides linear/non-linear solvers and time steppers

# libCEED Projects

Several projects built using libCEED

- Ratel - solid mechanics FEM and MPM (PSAAP)

- HONEE - fluid dynamics FEM & differential filtering (PHASTA)

- MFEM - various applications, libCEED integrators (LLNL)

- Palace - Electromagnetics FEM with MFEM + libCEED (Amazon)

- RDycore - River dynamical core with PETSc + libCEED (SciDAC)

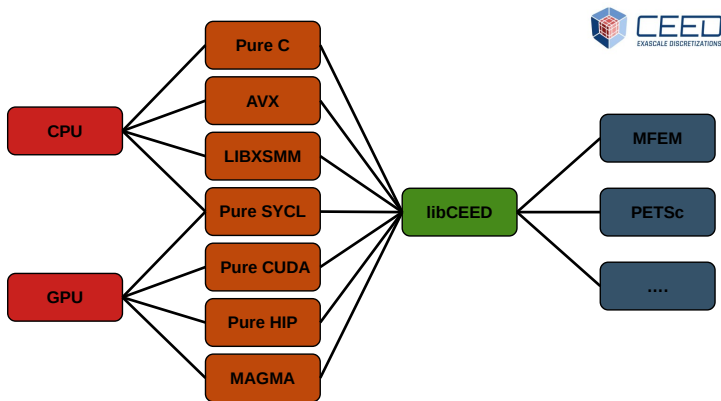# Matrix-Free Operators from libCEED

$$A = \mathcal{P}^T \mathcal{E}^T B^T D B \mathcal{E} \mathcal{P}$$



libCEED provides arbitrary order matrix-free operator evaluation

# Performance Portability from libCEED



libCEED backends target different hardware at runtime

# Extensible Solvers from PETSc

CeedEvaluator

MatCeed

CeedVector



libCEED provides the local operator action for PETSc objects

# Two Families of Approaches

Three libCEED backends with two approaches to operator application

- Separate kernels
  - /gpu/*/ref and /gpu/*/shared

  - $\mathcal{E}$, $B$, and $D$ all separate kernels

  - Higher overall memory usage, multiple kernel launches

- Fused kernel
  - /gpu/*/gen

  - Single kernel JiTed with data from $\mathcal{E}$, $B$, and $D$

  - Lower overall memory usage, single kernel launch

## Ref Operator Application

/gpu/*/ref and /gpu/*/shared use largely the same code

$A_L = \mathcal{E}^T B^T D B \mathcal{E}$ use separate kernels

$\mathcal{E}$ source comes from the /gpu/*/ref

/gpu/*/ref uses basic kernels for $B$

/gpu/*/shared uses shared memory for $B$
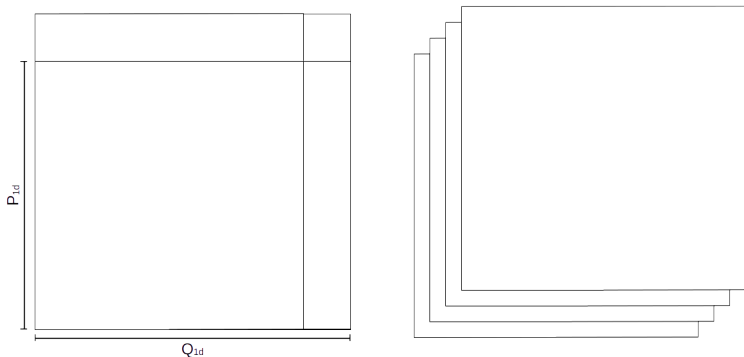
$D$ source is given by the user

## Shared Basis Code

```
 1  extern "C" __launch_bounds__(BASIS_GRAD_BLOCK_SIZE) __global__
 2      void Grad(const CeedInt num_elem, const CeedScalar *c_B,
 3      const CeedScalar *c_G, const CeedScalar *__restrict__ d_U,
 4      CeedScalar *__restrict__ d_V) {
 5    // Setup (omitted)
 6    // Apply basis element by element
 7    for (CeedInt elem=blockIdx.x*blockDim.z+threadIdx.z; elem<num_elem; elem+=gridDim.x
          *blockDim.z) {
 8      ReadElementStrided2d<NUM_COMP,P_1D>(data,elem,1,P_1D*P_1D*num_elem,P_1D*P_1D,d_U,
          r_U);
 9      GradTensor2d<NUM_COMP,P_1D,Q_1D>(data,r_U,s_B,s_G,r_V);
10      WriteElementStrided2d<NUM_COMP*DIM,Q_1D>(data,elem,1,Q_1D*Q_1D*num_elem,Q_1D*Q_1D
          ,r_V,d_V);
11    }
12  }
```

$x$ and $y$ thread index gives point (2D) or column of points (3D)

$z$ thread index gives the element

# Thread Usage



$x$ and $y$ thread index gives point (2D) or column of points (3D)

3D strategy works on 2D slabs of points

## Shared Basis Code

```
 1   template <int NUM_COMP, int P_1D, int Q_1D>
 2   inline __device__ void GradTensor2d(SharedData_Hip &data,
 3       const CeedScalar *__restrict__ r_U, const CeedScalar *c_B,
 4       const CeedScalar *c_G, CeedScalar *__restrict__ r_V) {
 5     CeedScalar r_t[1];
 6     for (CeedInt comp = 0; comp < NUM_COMP; comp++) {
 7       ContractX2d<NUM_COMP,P_1D,Q_1D>(data, &r_U[comp], c_G, r_t);
 8       ContractY2d<NUM_COMP,P_1D,Q_1D>(data, r_t, c_B, &r_V[comp+0*NUM_COMP]);
 9       ContractX2d<NUM_COMP,P_1D,Q_1D>(data, &r_U[comp], c_B, r_t);
10       ContractY2d<NUM_COMP,P_1D,Q_1D>(data, r_t, c_G, &r_V[comp+1*NUM_COMP]);
11     }
12   }
13
14   template <int NUM_COMP, int P_1D, int Q_1D>
15   inline __device__ void GradTransposeTensor2d(SharedData_Hip &data,
16       const CeedScalar *__restrict__ r_U, const CeedScalar *c_B,
17       const CeedScalar *c_G, CeedScalar *__restrict__ r_V) {
18     CeedScalar r_t[1];
19     for (CeedInt comp = 0; comp < NUM_COMP; comp++) {
20       ContractTransposeY2d<NUM_COMP,P_1D,Q_1D>(data, &r_U[comp+0*NUM_COMP], c_B, r_t);
21       ContractTransposeX2d<NUM_COMP,P_1D,Q_1D>(data, r_t, c_G, &r_V[comp]);
22       ContractTransposeY2d<NUM_COMP,P_1D,Q_1D>(data, &r_U[comp+1*NUM_COMP], c_G, r_t);
23       ContractTransposeAddX2d<NUM_COMP,P_1D,Q_1D>(data, r_t, c_B, &r_V[comp]);
24     }
25   }
```

Loop over components to reduce total shared memory needed

## Shared Basis Code

```
1  template <int NUM_COMP, int P_1D, int Q_1D>
2  inline __device__ void ContractX3d(SharedData_Hip &data, const CeedScalar *U,
3      const CeedScalar *B, CeedScalar *V) {
4    CeedScalar r_B[P_1D];
5    for (CeedInt i = 0; i < P_1D; i++) r_B[i] = B[i + data.t_id_x * P_1D];
6
7    for (CeedInt k = 0; k < P_1D; k++) {
8      data.slice[data.t_id_x + data.t_id_y * T_1D] = U[k];
9      __syncthreads();
10     V[k] = 0.0;
11     if (data.t_id_x < Q_1D && data.t_id_y < P_1D) {
12       for (CeedInt i = 0; i < P_1D; i++) {
13         V[k] += r_B[i] * data.slice[i + data.t_id_y * T_1D];
14       }
15     }
16     __syncthreads();
17   }
18 }
```

Each thread computes all node's contributions to one quadrature point

3D loops over 2D slabs for tensor contraction

# Gen Operator Application

/gpu/*/gen generates a single kernel for the operator

$A_L = \mathcal{E}^T B^T D B \mathcal{E}$ uses a single kernel

$\mathcal{E}$ source comes from the /gpu/*/ref

$B$ source comes from /gpu/*/shared

$D$ source is given by the user

## Generated Operator Kernel

```
1  extern "C" __global__ void CeedKernelCudaGenOperator_mass (CeedInt num_elem,
2      void* ctx, FieldsInt_Cuda indices, Fields_Cuda fields, Fields_Cuda B,
3      Fields_Cuda G, CeedScalar *W, Points_Cuda points) {
4    // Setup kernel data
5
6    // Input and Output field constants and basis data
7
8    // Element loop
9    __syncthreads();
10   for (CeedInt elem = blockIdx.x*blockDim.z + threadIdx.z; elem < num_elem; elem +=
          gridDim.x*blockDim.z) {
11     // -- Input field restrictions (E) and basis actions (B)
12
13     // -- Output field setup
14     {
15       // -- Apply QFunction (D)
16       mass(ctx, 1, inputs, outputs);
17     }
18
19     // -- Output field basis actions (B^T) and restrictions (E^T)
20   }
21 }
22 // ----------------------------------------------------------------------------
```

$$A_L = \mathcal{E}^T B^T D B \mathcal{E} \text{ in a single kernel}$$

## Generated Operator Kernel

```
1    // Setup kernel data
2    const CeedScalar *d_in_0 = fields.inputs[0];
3    const CeedScalar *d_in_1 = fields.inputs[1];
4    CeedScalar *d_out_0 = fields.outputs[0];
5    const CeedInt dim = 1;
6    const CeedInt Q_1d = 8;
7    extern __shared__ CeedScalar slice[];
8    SharedData_Cuda data;
9    data.t_id_x = threadIdx.x;
10   data.t_id_y = threadIdx.y;
11   data.t_id_z = threadIdx.z;
12   data.t_id  = threadIdx.x+threadIdx.y*blockDim.x+threadIdx.z*blockDim.y*blockDim.x;
13   data.slice = slice + data.t_id_z*T_1D;
```

Set up pointers to basis data and shared memory

# Generated Operator Kernel

```
1     // Input field constants and basis data
2     // -- Input field 0
3     const CeedInt P_1d_in_0 = 8;
4     const CeedInt num_comp_in_0 = 1;
5     // EvalMode: none
6     // -- Input field 1
7     const CeedInt P_1d_in_1 = 5;
8     const CeedInt num_comp_in_1 = 1;
9     // EvalMode: interpolation
10    __shared__ CeedScalar s_B_in_1[40];
11    LoadMatrix<P_1d_in_1, Q_1d>(data, B.inputs[1], s_B_in_1);
12
13    // Output field constants and basis data
14    // -- Output field 0
15    const CeedInt P_1d_out_0 = 5;
16    const CeedInt num_comp_out_0 = 1;
17    // EvalMode: interpolation
18    __shared__ CeedScalar s_B_out_0[40];
19    LoadMatrix<P_1d_out_0, Q_1d>(data, B.outputs[0], s_B_out_0);
```

Basis data and constants loaded

## Generated Operator Kernel

```
1        // Scratch restriction buffer space
2        CeedScalar r_e_scratch[8];
3
4        // -- Input field restrictions and basis actions
5        // ---- Input field 0
6        CeedScalar r_e_in_0[num_comp_in_0*P_1d_in_0];
7        // Strides: {1, 8, 8}
8        ReadLVecStrided1d<num_comp_in_0, P_1d_in_0,1,8,8>(data, elem, d_in_0, r_e_in_0);
9        // EvalMode: none
10       CeedScalar *r_q_in_0 = r_e_in_0;
11       // ---- Input field 1
12       CeedScalar *r_e_in_1 = r_e_scratch;
13       const CeedInt l_size_in_1 = 61;
14       // CompStride: 1
15       ReadLVecStandard1d<num_comp_in_1, 1, P_1d_in_1>(data, l_size_in_1, elem, indices.
             inputs[1], d_in_1, r_e_in_1);
16       // EvalMode: interpolation
17       CeedScalar r_q_in_1[num_comp_in_1*Q_1d];
18       Interp1d<num_comp_in_1, P_1d_in_1, Q_1d>(data, r_e_in_1, s_B_in_1, r_q_in_1);
19
20       // -- Output field setup
21       // ---- Output field 0
22       CeedScalar r_q_out_0[num_comp_out_0*Q_1d];
```

Restrict and apply basis for each input

Setup output data buffers

## Generated Operator Kernel

```
1         // Note: Using full elements
2         {
3           // -- Input fields
4           // ---- Input field 0
5           CeedScalar *r_s_in_0 = r_q_in_0;
6           // ---- Input field 1
7           CeedScalar *r_s_in_1 = r_q_in_1;
8           // -- Output fields
9           // ---- Output field 0
10          CeedScalar *r_s_out_0 = r_q_out_0;
11
12          // -- QFunction inputs and outputs
13          // ---- Inputs
14          CeedScalar *inputs[2];
15          // ------ Input field 0
16          inputs[0] = r_s_in_0;
17          // ------ Input field 1
18          inputs[1] = r_s_in_1;
19          // ---- Outputs
20          CeedScalar *outputs[1];
21          // ------ Output field 0
22          outputs[0] = r_s_out_0;
23
24          // -- Apply QFunction
25          mass(ctx, 1, inputs, outputs);
26        }
```
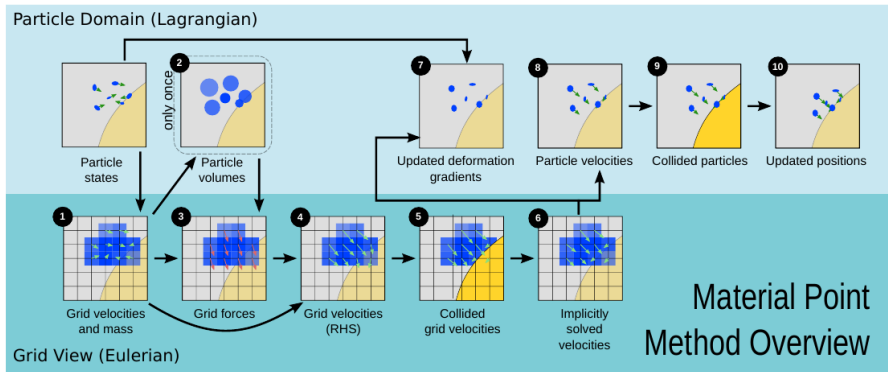
Apply QFunction at each quadrature point
May apply to 2D slabs or full elements in 3D

# Generated Operator Kernel

```
1        // -- Output field basis action and restrictions
2        // ---- Output field 0
3        // EvalMode: interpolation
4        CeedScalar *r_e_out_0 = r_e_scratch;
5        InterpTranspose1d<num_comp_out_0, P_1d_out_0, Q_1d>(data, r_q_out_0, s_B_out_0,
               r_e_out_0);
6        const CeedInt l_size_out_0 = 61;
7        // CompStride: 1
8        WriteLVecStandard1d<num_comp_out_0, 1, P_1d_out_0>(data, l_size_out_0, elem,
               indices.outputs[0], r_e_out_0, d_out_0);
```

Output basis action and restriction to assemble result

# What is MPM?



Particle Domain (Lagrangian)

only once

Particle states

Particle volumes

Updated deformation gradients

Particle velocities

Collided particles

Updated positions

Grid velocities and mass

Grid forces

Grid velocities (RHS)

Collided grid velocities

Implicitly solved velocities

Grid View (Eulerian)
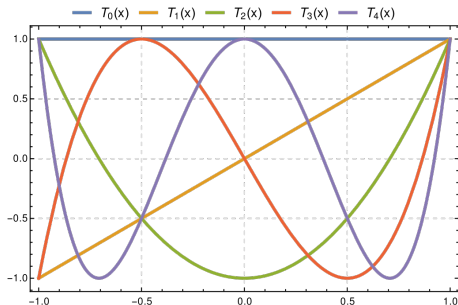
Material Point Method Overview

- Continuum based particle method with background mesh for gradients
- Extension of FLIP (which is an extension of PIC)
- Used in rendering for the movie *Frozen*

# What does MPM have to do with FEM?

- Problem on background mesh changes when material points move

- Natural fit for matrix-free representation

- Similar reasoning to use matrix-free for adaptive methods

- Ratel/libCEED FEM infrastructure provides fast background mesh solves

# libCEED Basis Evaluation AtPoints



- Interpolate from primal to dual (quadrature) space
- Fit Chebyshev polynomials to values at quadrature points
- Evaluate Chebyshev polynomials at reference coords of material points
- Transpose the order for projection to mesh from material points

## Shared Basis Code

```
1  extern "C" __launch_bounds__(BASIS_INTERP_BLOCK_SIZE) __global__
2     void InterpAtPoints(const CeedInt num_elem, const CeedScalar *c_B,
3     const CeedInt *points_per_elem, const CeedScalar *__restrict__ d_X,
4     const CeedScalar *__restrict__ d_U, CeedScalar *__restrict__ d_V) {
5    // Setup (omitted)
6    // Apply basis element by element
7    for (CeedInt elem=blockIdx.x*blockDim.z+threadIdx.z; elem<num_elem; elem+=gridDim.x
         *blockDim.z) {
8      // Map from nodes to Chebyshev coefficients
9      ReadElementStrided2d<NUM_COMP,P_1D>(data,elem,1,P_1D*P_1D*num_elem,P_1D*P_1D,d_U,
           r_U);
10     InterpTensor2d<NUM_COMP,P_1D,Q_1D>(data,r_U,s_B,r_C);
11     // Map from Chebyshev coefficients to points
12     for (CeedInt i=threadIdx.x+threadIdx.y*blockDim.x; i<point_loop_bound; i+=
           blockDim.x*blockDim.y) {
13       const CeedInt p = i % NUM_PTS;
14       ReadPoint<DIM,NUM_PTS>(data,elem,p,NUM_PTS,1,num_elem*NUM_PTS,NUM_PTS,d_X,r_X);
15       InterpAtPoints2d<NUM_COMP,NUM_PTS,Q_1D>(data,i,r_C,r_X,r_V);
16       WritePoint<NUM_COMP,NUM_PTS>(data,elem,p,NUM_PTS,1,num_elem*NUM_PTS,NUM_PTS,r_V
             ,d_V);
17     }
18   }
19 }
```

Threadblock maps to Chebeshev coeffs on element (standard interpolation)

Each thread maps from Chebyshev coeffs to single point

# Shared Basis Code

```
1   template <int NUM_COMP, int NUM_POINTS, int Q_1D>
2   inline __device__ void InterpAtPoints2d(SharedData_Hip &data, const CeedInt p, const
        CeedScalar *__restrict__ r_C, const CeedScalar *__restrict__ r_X,
3                                           CeedScalar *__restrict__ r_V) {
4     for (CeedInt i = 0; i < NUM_COMP; i++) r_V[i] = 0.0;
5     for (CeedInt comp = 0; comp < NUM_COMP; comp++) {
6       CeedScalar buffer[Q_1D];
7       CeedScalar chebyshev_x[Q_1D];
8       // Load coefficients
9       if (data.t_id_x<Q_1D && data.t_id_y<Q_1D) {
10        data.slice[data.t_id_x+data.t_id_y*Q_1D] = r_C[comp];
11      }
12      __syncthreads();
13      // Contract x direction
14      ChebyshevPolynomialsAtPoint<Q_1D>(r_X[0], chebyshev_x);
15      for (CeedInt i = 0; i < Q_1D; i++) {
16        buffer[i] = 0.0;
17        for (CeedInt j = 0; j < Q_1D; j++) {
18          buffer[i] += chebyshev_x[j] * data.slice[j + i*Q_1D];
19        }
20      }
21      // Contract y direction
22      ChebyshevPolynomialsAtPoint<Q_1D>(r_X[1], chebyshev_x);
23      for (CeedInt i = 0; i < Q_1D; i++) {
24        r_V[comp] += chebyshev_x[i] * buffer[i];
25      }
26    }
27  }
```

Each thread (point) needs separate contraction buffers

# Preconditioning Support

Some operator assembly needed for preconditioning

Diagonal assembly for Jacobi, Chebyshev

Full assembly for AMG or LU

# FEM Diagonal Assembly

FEM diagonal assembly consists of two phases

- QFunction assembly
  - Assemble small matrix at each quadrature point
  - Each active input individually set to 1
  - QFunction kernel $D$ called to populate assembled row

- Diagonal assembly
  - Compute diagonal entries of $B^T D B$ on element
  - Element diagonals assembled via $\mathcal{E}^T$ for local diagonal

# MPM Diagonal Assembly
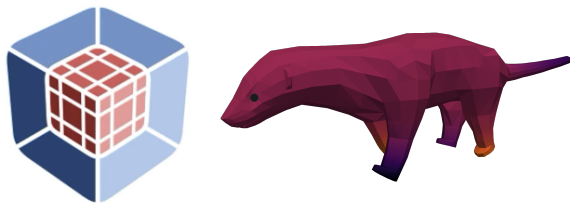
MPM diagonal assembly consists of a single phase

- Each input node on elements individually set to 1

- Basis $B$ applied to get values at material points

- QFunction kernel $D$ called to populate result

- Basis $B^T$ applied to get values at nodes

- Corresponding diagonal entry copied into element diagonal

- Element diagonals assembled via $\mathcal{E}^T$ for local diagonal

# FEM Full Assembly

FEM full assembly consists of three phases

- Sparsity pattern
  - Compute local (on process) matrix COO indices for entries

- QFunction assembly
  - Same as for diagonal assembly

- Full assembly
  - Compute entries of $B^T D B$ for each element
  - Populate assembled array per the sparsity pattern

## Questions?



libCEED Repo: https://github.com/CEED/libCEED
HONEE Repr: https://gitlab.com/phypid/HONEE
Ratel Repo: https://gitlab.com/micromorph/ratel