

# Empowering MFEM Using libCEED: Features and Performance Analysis

MFEM Community Workshop – Livermore, California, USA

Yohann Dudouit<sup>1</sup>, Natalie Beams<sup>2</sup>, Jed Brown<sup>3</sup>, John Camier<sup>1</sup>, Veselin Dobrev<sup>1</sup>, Tzanio Kolev<sup>1</sup>, Jeremy Thompson<sup>3</sup>, Tim Warburton<sup>4</sup>, & the CEED Team<sup>1,2,3,4,5,6,7</sup>

October 26, 2023



# The CEED project: The partial assembly decomposition

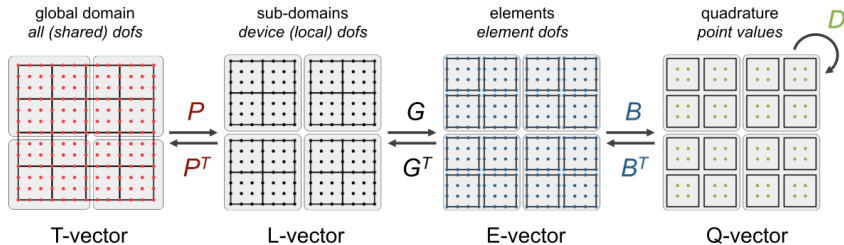


CEED  
EXASCALE DISCRETIZATIONS

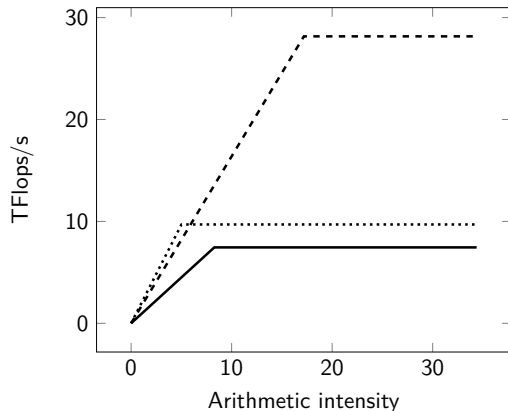
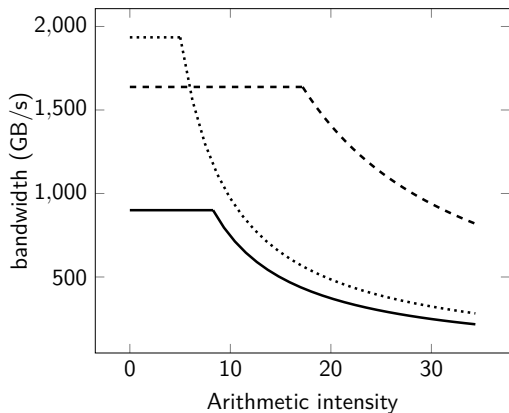


The assembly/evaluation of FEM operators can be decomposed into **parallel**, **mesh topology**, **basis**, and **geometry/physics** components:

$$A = P^T G^T B^T D B G P$$

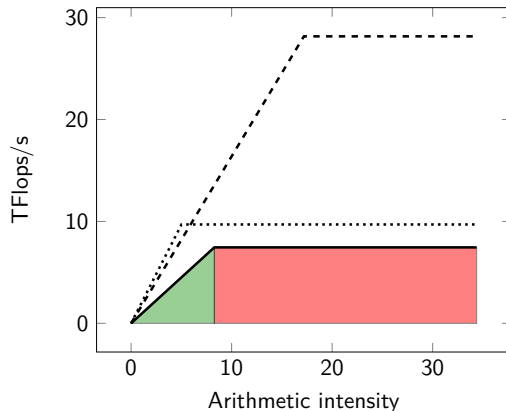
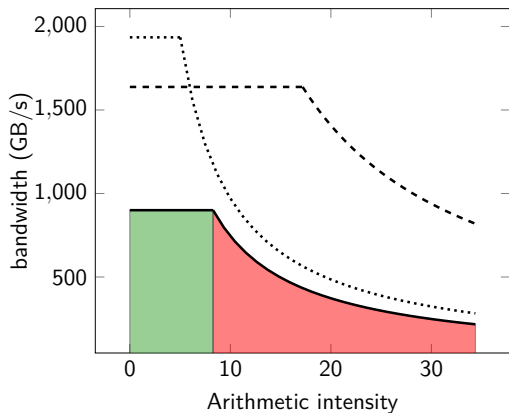


# Roofline model: The two sides of GPU performance



- NVIDIA V100
- ..... NVIDIA A100
- - - AMD MI250X

# Roofline model: The two sides of GPU performance

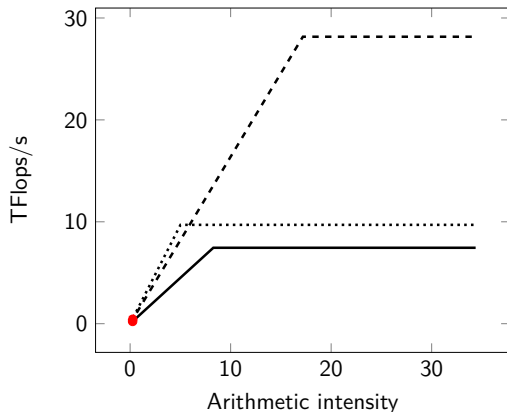
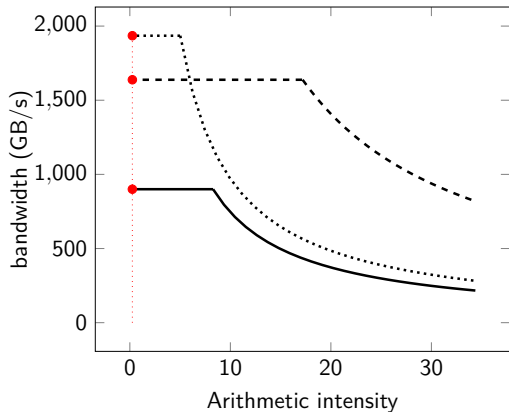


- NVIDIA V100
- ..... NVIDIA A100
- - - AMD MI250X
- Memory bound
- Computation bound

$$t_M = \frac{\text{Data}}{\text{MaxBW}}$$

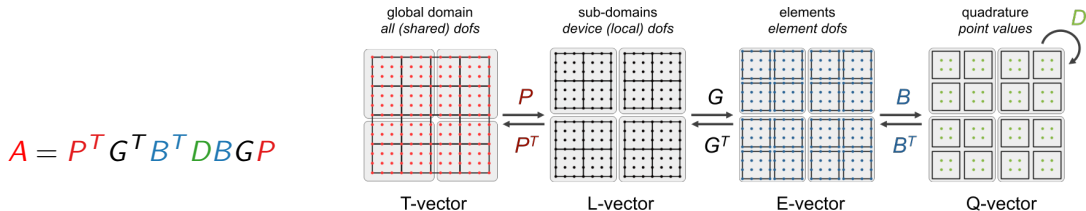
$$t_C = \frac{\text{Flops}}{\text{MaxFlops}}$$

## Roofline model: The sparse matrix case



⇒ The only way to run faster than a sparse matrix is to **move less data**.

# Partial assembly: The algorithmic costs



	Sparse Matrix	Partial Assembly	$G$	$B$	$D$
Amount of storage	$\mathcal{O}(np^{2d})$	$\mathcal{O}(np^d)$	$\mathcal{O}(np^d)$	$\mathcal{O}(p^{2d})$	$\mathcal{O}(np^d)$
FLOPs to apply	$\mathcal{O}(np^{2d})$	$\mathcal{O}(np^{2d})$	$\mathcal{O}(np^d)$	$\mathcal{O}(np^{2d})$	$\mathcal{O}(np^d)$
Arithmetic intensity	$\mathcal{O}(1)$	$\mathcal{O}(p^d)$	$\mathcal{O}(1)$	$\mathcal{O}(p^d)$	$\mathcal{O}(1)$

## Potential speedup:

Data movement and storage is reduced from  $\mathcal{O}(np^{2d})$  to  $\mathcal{O}(np^d)$  to apply the finite element operator, potential speedup for Partial Assembly of  $\sim \mathcal{O}(p^d)$ .

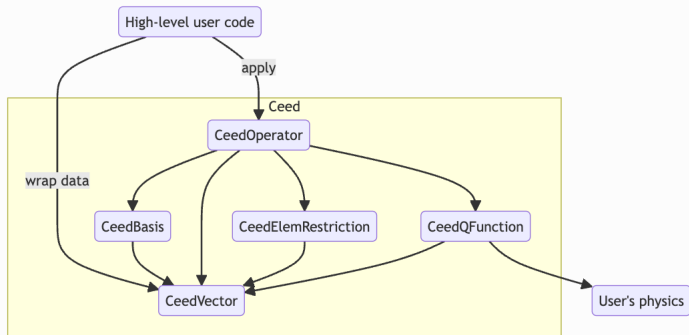
## Reducing the arithmetic intensity: The sum factorization trick

On tensor product finite elements, the  $B$  operator can be computed as:

$$\begin{aligned}v_{k_1 k_2 k_3} &= B_{IK} u_I = \underbrace{\sum_{i_1, i_2, i_3} u_{i_1 i_2 i_3} \varphi_{i_1}(x_{k_1}) \varphi_{i_2}(x_{k_2}) \varphi_{i_3}(x_{k_3})}_{\mathcal{O}(p^6) = \mathcal{O}(p^{2d})} \\ &= \sum_{i_3} \varphi_{i_3}(x_{k_3}) \left( \sum_{i_2} \varphi_{i_2}(x_{k_2}) \left( \underbrace{\sum_{i_1} \varphi_{i_1}(x_{k_1}) u_{i_1 i_2 i_3}}_{\mathcal{O}(p^4) = \mathcal{O}(p^{d+1})} \right) \right) \\ &= \tilde{B}_{i_3 k_3} \otimes \tilde{B}_{i_2 k_2} \otimes \tilde{B}_{i_1 k_1} u_{i_1 i_2 i_3}\end{aligned}$$

	No Sum Factorization	Sum Factorization
Amount of storage for $B$	$\mathcal{O}(p^{2d})$	$\mathcal{O}(p^2)$
FLOPs to apply	$\mathcal{O}(p^{2d})$	$\mathcal{O}(p^{d+1})$
Arithmetic intensity	$\mathcal{O}(p^d)$	$\mathcal{O}(p)$

# The libCEED core interface



CeedOperator	$A = P^T G^T B^T D B G P$
CeedBasis	$B$
CeedElemRestriction	$G$
CeedQFunction	$D$
CeedVector	Wrapper for degrees-of-freedom/data at quadrature points



# Features of the libCEED library

Backend	Description
AVX	Optimized cpu backend taking advantage of AVX instructions.
CUDA	Pure CUDA backend using JIT compilation.
HIP	Pure HIP backend using JIT compilation.
SYCL	Pure SYCL backend using JIT compilation.
Magma	Backend leveraging the Magma library, high performance on non-tensor elements.
XSMW	Backend leveraging the libXSMW library, highest cpu performance.
OCCA	Backend based on the OCCA abstraction layer.

## Extra libCEED features:

- Provide an interface to compute the **diagonal** of any operator,
- Provide an interface to **assemble a sparse-matrix** for any operator,
- Provide an interface for p-multigrid (Jeremy Thompson).

# Features of the libCEED integration in MFEM

## Using the libCEED backend:

- **MFEM\_USE\_CCEED=YES.**
- **-d ceed-cpu/ceed-cuda/ceed-hip.**
- Specific libCEED backends can be selected using :, e.g.  
**-d ceed-hip:/gpu/hip/magma.**

Supported MFEM Integrators	Weak form
MassIntegrator	$\int u v$
VectorMassIntegrator	$\int \mathbf{u} \cdot \mathbf{v}$
ConvectionIntegrator	$\int (\mathbf{a} \cdot \nabla u) v$
VectorConvectionNLFIntegrator	$\int c(\nabla \mathbf{u}) \cdot \mathbf{v}$
DiffusionIntegrator	$\int c \nabla u \cdot \nabla v$
VectorDiffusionIntegrator	$\int c \nabla \mathbf{u} \cdot \nabla \mathbf{v}$

## Pros of the libCEED backend:

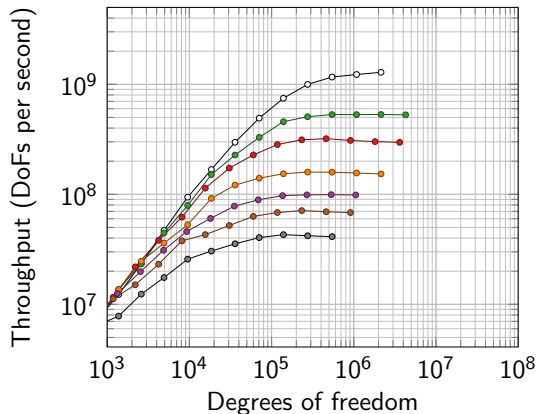
- Support for mixed-meshes (including simplices) and p-adaptivity (limited to serial),
- Interface to construct partial assembly and fully matrix-free operators,
- Algebraic multigrid solver based on the libCEED interface (Andrew Barker).

## Cons of the libCEED backend:

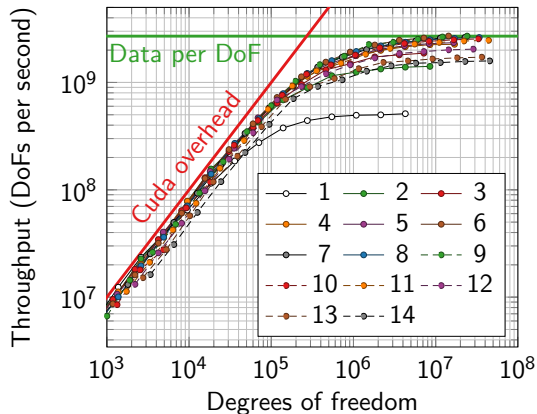
- Does not currently support as many integrators as native MFEM,
- libCEED GPU operators can be "non-deterministic" (use atomic operations).

# Comparing sparse-matrix and matrix-free

Sparse Matrix-Vector Product ( CuSparse )

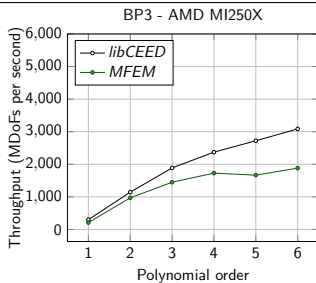
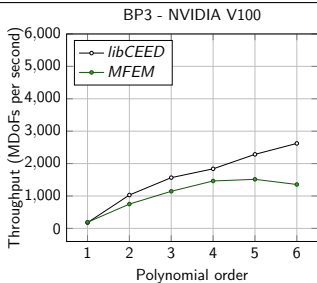
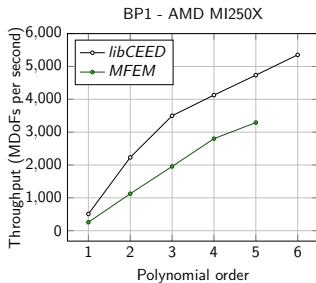
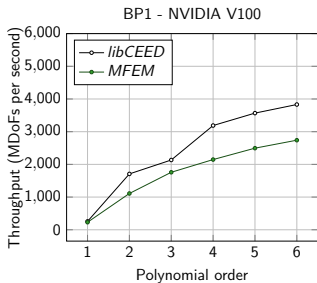


Matrix-Free using libCEED (cuda-gen)

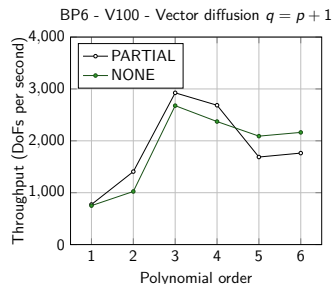
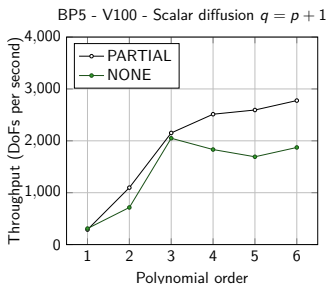
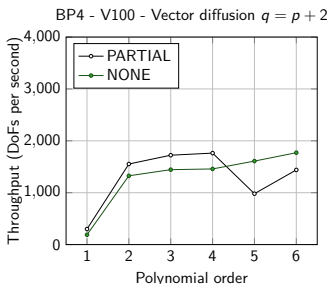
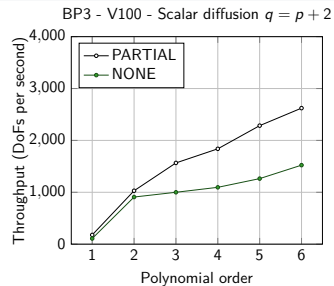
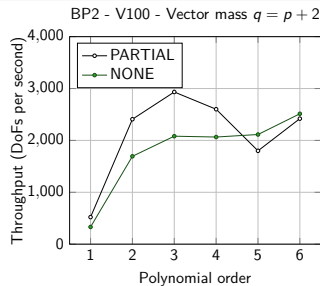
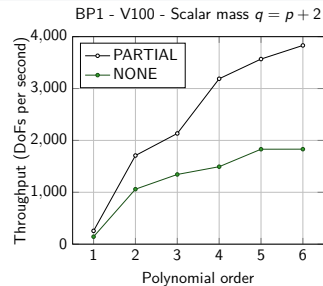


$$t_M = \frac{\text{Data}}{\text{MaxBW}} \Rightarrow \text{MaxThroughput} = \frac{\text{MaxBW}}{\text{Data per DoF}}$$

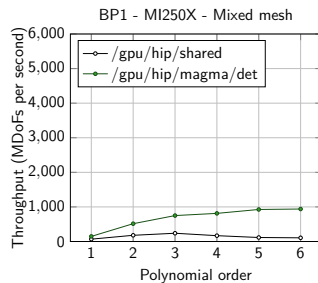
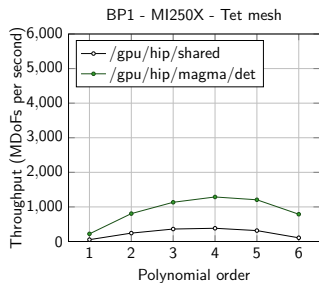
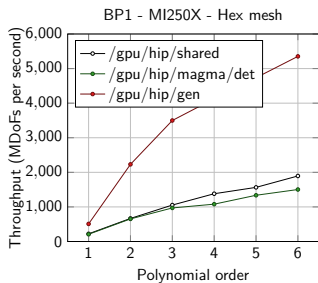
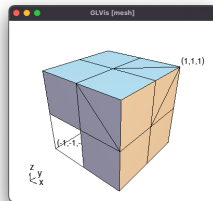
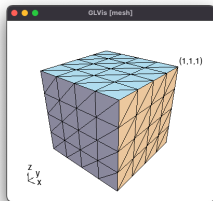
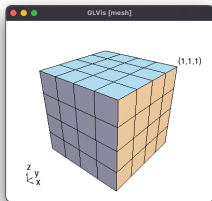
# Comparing native MFEM and the libCEED backend



# Comparing the assembly levels: AssemblyLevel::PARTIAL vs AssemblyLevel::NONE



# Performance on simplicies and mixed-meshes



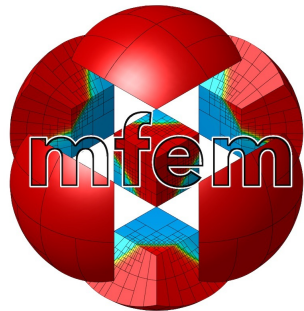
# Key features and future directions

## Key features:

- Competitive performance,
- Run on any hardware ( cpu, CUDA, HIP, SYCL ),
- Support for simplices and mixed-meshes,
- Support for p-adaptivity.

## Future directions:

- Add support for  $H(\text{div})$  and  $H(\text{curl})$  (non-tensor only),
- Add support for discontinuous Galerkin methods,
- Add support for sparse-matrix assembly through libCEED.



**CEED**  
EXASCALE DISCRETIZATIONS



# Disclaimer

Disclaimer: This document was prepared as an account of work sponsored by an agency of the United States government, Neither the United States government or Lawrence Livermore National Security, LLC, nor any of their employees make 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 represent 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.