

On performance portability and library reuse: A Navier-Stokes miniapp

Valeria Barra¹, Jed Brown¹, Jeremy Thompson¹, Yohann Dudouit²

¹ University of Colorado at Boulder

² Lawrence Livermore National Lab

University of Colorado at Boulder (Computer Science Dept.)

ICIAM 2019, Valencia, Spain — July 17, 2019



CENTER FOR EFFICIENT
EXASCALE DISCRETIZATIONS



Overview

- We will introduce libCEED: the library of the CEED
- For decades, high-order numerical methods have been considered too expensive. But we now know that a sparse matrix is no longer a good representation for high-order operators
- libCEED uses a matrix-free operator description, based on a purely algebraic interface, where user only specifies action of weak form operators
- libCEED operator representation is optimal with respect to the FLOPs needed for its evaluation, as well as the memory transfer needed for operator evaluations (matvec)
 - For a 3D Jacobian operator, matrix-free operators that exploit tensor-product structures reduce the work load from $O(p^6)$ (for sparse matrix) to $O(p^4)$, and memory storage from $O(p^6)$ to $O(p^3)$
- We demonstrate the usage of libCEED with PETSc for a compressible Navier-Stokes solver



libCEED: the library of the CEED

(Center for Efficient Extensible Discretizations)

- Primary target: high-order finite element/spectral element methods exploiting tensor product structure
- Open source (BSD-2 license) C library with Fortran interface
- Releases: v0.1 (January 2018), v0.2 (March 2018), v0.3 (September 2018), v0.4 (March 2019)

For latest release:

Brown J., Abdelfattah A., **Barra V.**, Dobrev V., Dudouit Y., Fischer P., et al., *CEED ECP Milestone Report: Public release of CEED 2.0* (2019, March 31) DOI:
<http://doi.org/10.5281/zenodo.2641316>



University of Colorado
Boulder

libCEED: the library of the CEED (Center for Efficient Extensible Discretizations)

<code>/cpu/self/ref/*:</code>	with * reference serial and blocked implementations
<code>/cpu/self/avx/*:</code>	AVX (Advanced Vector Extensions instruction sets) with * reference serial and blocked implementations
<code>/cpu/self/xsmm/*:</code>	LIBXSMM (Intel library for small dense/sparse mat-multiply) with * reference serial and blocked implementations
<code>/*/occa:</code>	OCCA (just-in-time compilation) with *: CPU, GPU, OpenMP (Open Multi-Processing: API), OpenCL (framework for CPUs, GPUs, etc.)
<code>/gpu/magma:</code>	MAGMA (dense Linear Algebra library for GPUs and multicore architectures)
<code>/gpu/cuda/*:</code>	CUDA with *: <code>ref</code> (reference pure CUDA kernels), <code>reg</code> (pure CUDA kernels using one thread per element), <code>shared</code> , optimized pure CUDA kernels using shared memory

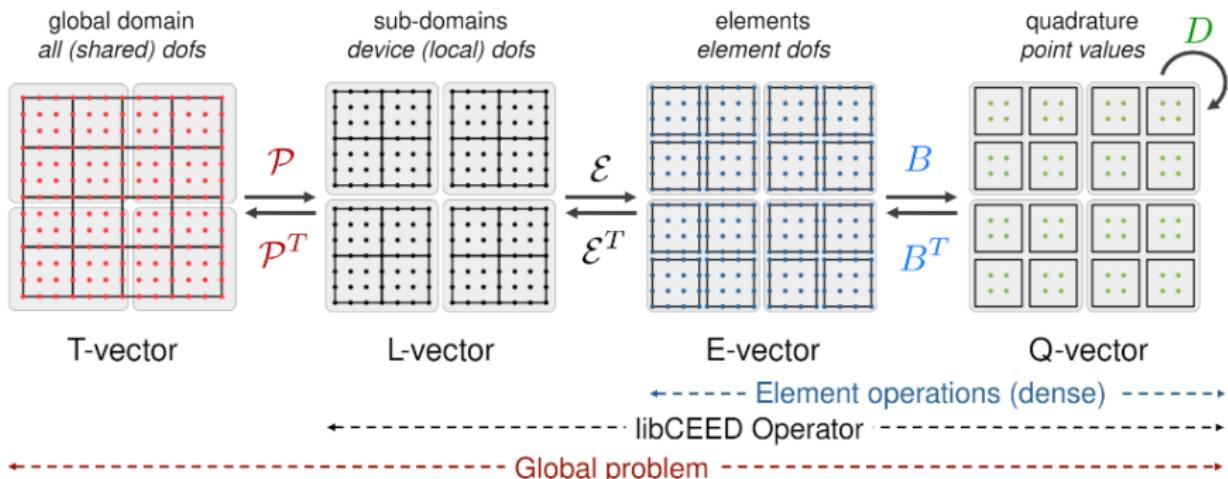
Same source code can call multiple CEEDs with different backends. On-device operator implementation with unique interface



libCEED decomposition



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



Performance w. r. t. space, on KNL

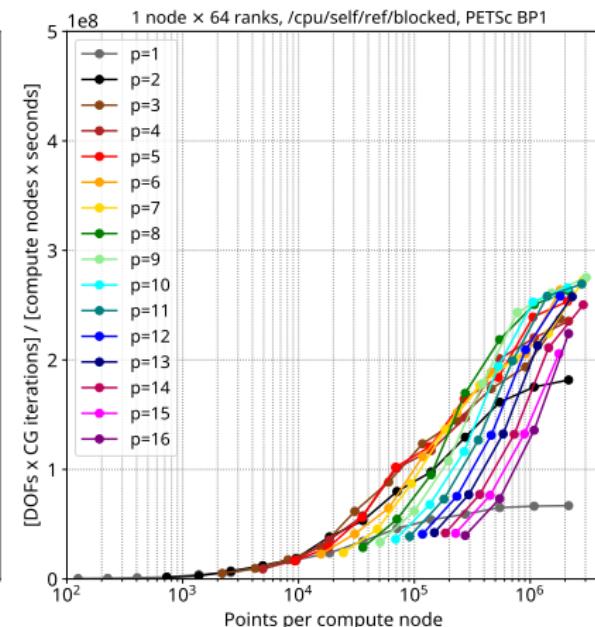
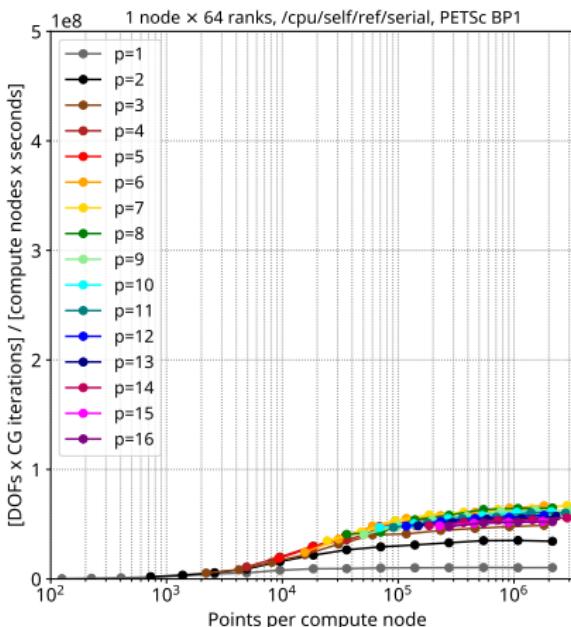
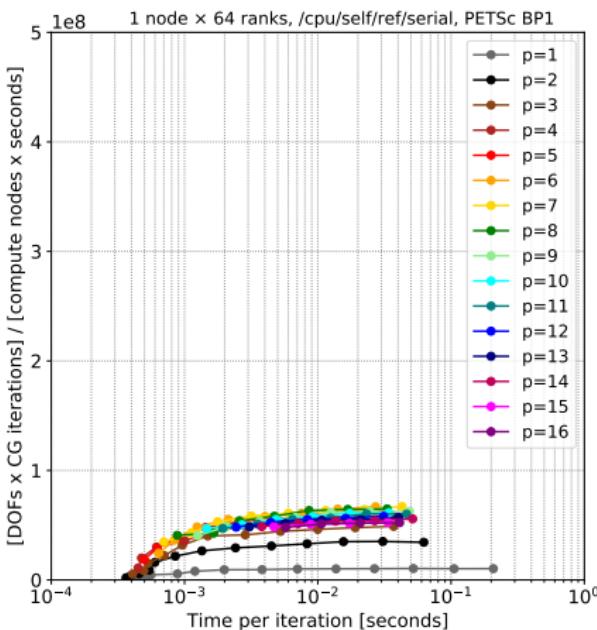


Figure: Knight Landing (Intel Xeon Phi 7230 SKU 1.3 GHz) with intel-18 compiler. In (a) serial implementation; in (b) blocked implementation ($q = P + 1$, $P = p + 1$)

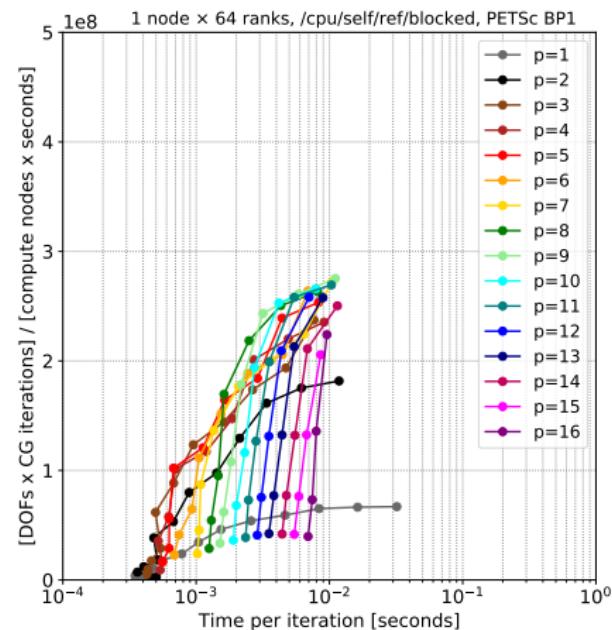


University of Colorado
Boulder

Performance w. r. t. time, on KNL



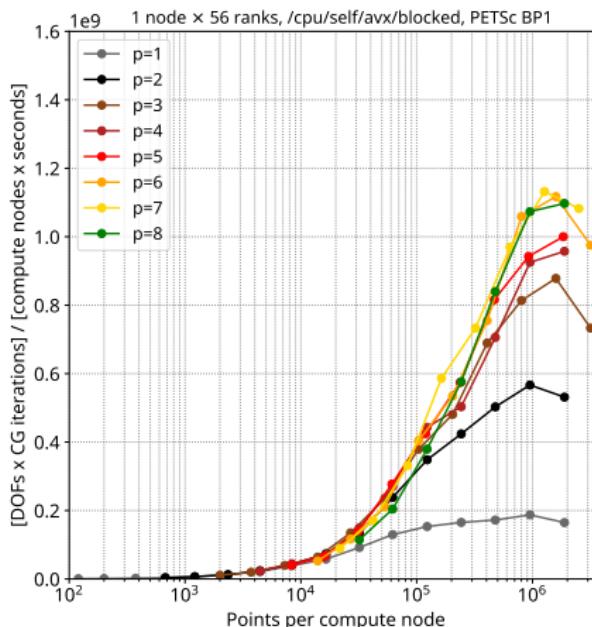
(a)



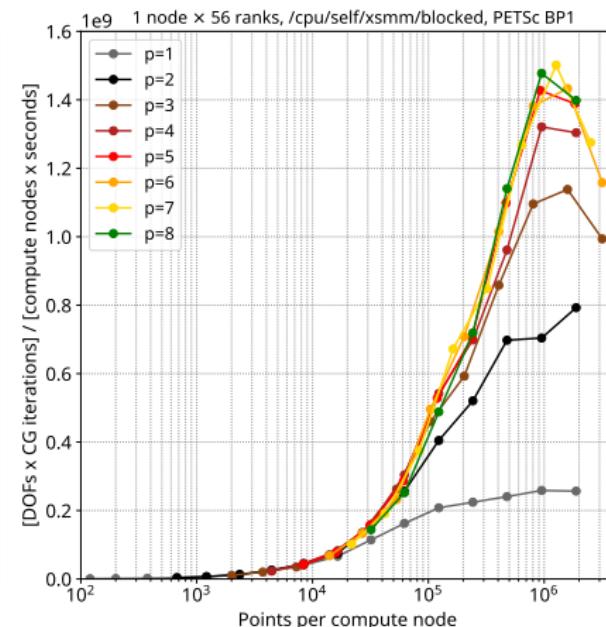
(b)

Figure: Knight Landing (Intel Xeon Phi 7230 SKU 1.3 GHz) with intel-18 compiler. In (a) serial implementation; in (b) blocked implementation ($q = P + 1$, $P = p + 1$)

Performance w. r. t. space, on Skylake: blocked



(a)



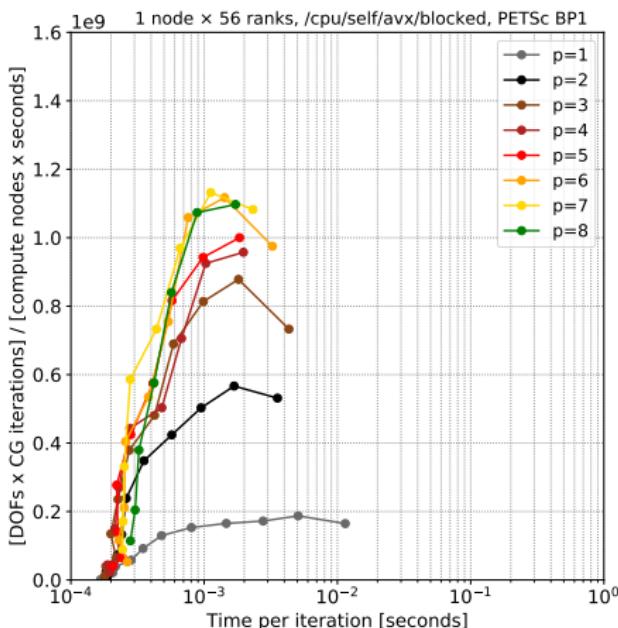
(b)

Figure: Skylake (2x Intel Xeon Platinum 8180M CPU 2.50GHz) with intel-19 compiler.
Backends: in (a) AVX; in (b) libXSMM ($q = P + 1$, $P = p + 1$)

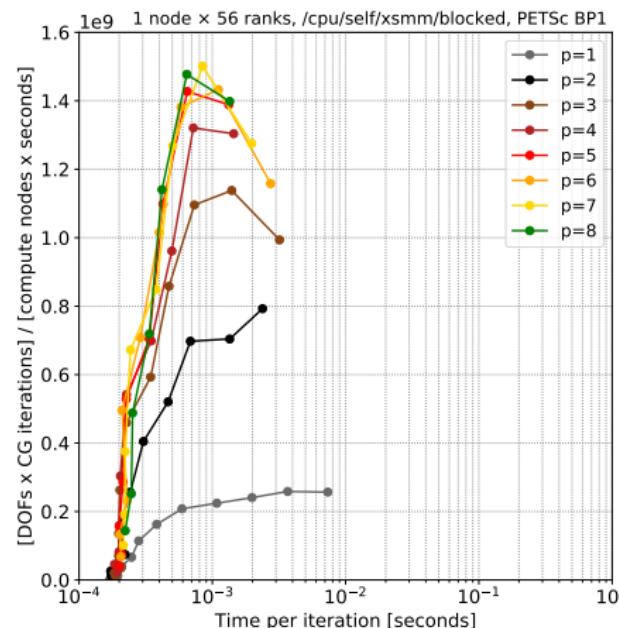


University of Colorado
Boulder

Performance w. r. t. time, on Skylake: blocked



(a)



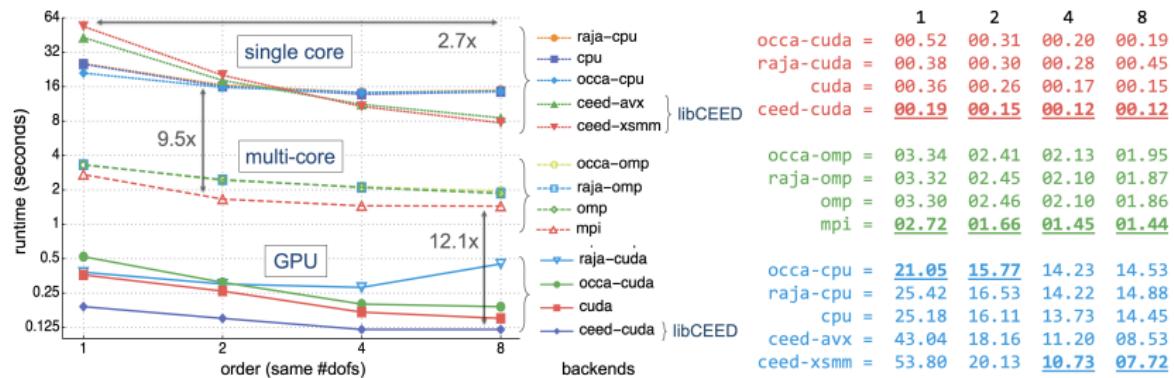
(b)

Figure: Skylake (2x Intel Xeon Platinum 8180M CPU 2.50GHz) with intel-19 compiler.
Backends: in (a) AVX; in (b) libXSMM ($q = P + 1$, $P = p + 1$)



University of Colorado
Boulder

Preliminary GPU results: MFEM + libCEED



single-GPU, multi-core CPU, and single-core CPU,
for 1.3 millions DOFs in 2D.

Results by Yohann Dudouit on a Linux desktop with a Quadro GV100 GPU,
sm_70, CUDA 10.1, and Intel Xeon Gold 6130 CPU @ 2.10GHz.



Towards a libCEED miniapp: a Navier-Stokes solver

Compressible Navier-Stokes equations in conservation form:

$$\frac{\partial \rho}{\partial t} + \nabla \cdot \mathbf{u} = 0, \quad (1a)$$

$$\frac{\partial \mathbf{u}}{\partial t} + \nabla \cdot \left(\frac{\mathbf{u} \otimes \mathbf{u}}{\rho} + P \mathbf{I}_3 \right) + \rho g \mathbf{k} = \nabla \cdot \boldsymbol{\sigma}, \quad (1b)$$

$$\frac{\partial E}{\partial t} + \nabla \cdot \left(\frac{(E + P)\mathbf{u}}{\rho} \right) = \nabla \cdot (\mathbf{u} \cdot \boldsymbol{\sigma} + k \nabla T), \quad (1c)$$



Towards a libCEED miniapp: a Navier-Stokes solver

Compressible Navier-Stokes equations in conservation form:

$$\frac{\partial \rho}{\partial t} + \nabla \cdot \mathbf{u} = 0, \quad (1a)$$

$$\frac{\partial \mathbf{u}}{\partial t} + \nabla \cdot \left(\frac{\mathbf{u} \otimes \mathbf{u}}{\rho} + P \mathbf{I}_3 \right) + \rho g \mathbf{k} = \nabla \cdot \boldsymbol{\sigma}, \quad (1b)$$

$$\frac{\partial E}{\partial t} + \nabla \cdot \left(\frac{(E + P)\mathbf{u}}{\rho} \right) = \nabla \cdot (\mathbf{u} \cdot \boldsymbol{\sigma} + k \nabla T), \quad (1c)$$

where $\boldsymbol{\sigma} = \mu(\nabla \mathbf{u} + (\nabla \mathbf{u})^T + \lambda(\nabla \cdot \mathbf{u})\mathbf{I}_3)$, and



Towards a libCEED miniapp: a Navier-Stokes solver

Compressible Navier-Stokes equations in conservation form:

$$\frac{\partial \rho}{\partial t} + \nabla \cdot \mathbf{u} = 0, \quad (1a)$$

$$\frac{\partial \mathbf{u}}{\partial t} + \nabla \cdot \left(\frac{\mathbf{u} \otimes \mathbf{u}}{\rho} + P \mathbf{I}_3 \right) + \rho g \mathbf{k} = \nabla \cdot \boldsymbol{\sigma}, \quad (1b)$$

$$\frac{\partial E}{\partial t} + \nabla \cdot \left(\frac{(E + P)\mathbf{u}}{\rho} \right) = \nabla \cdot (\mathbf{u} \cdot \boldsymbol{\sigma} + k \nabla T), \quad (1c)$$

where $\boldsymbol{\sigma} = \mu(\nabla \mathbf{u} + (\nabla \mathbf{u})^T + \lambda(\nabla \cdot \mathbf{u})\mathbf{I}_3)$, and

$$(c_p/c_v - 1)(E - \mathbf{u} \cdot \mathbf{u}/(2\rho) - \rho gz) = P \quad \leftarrow \text{pressure}$$

μ \leftarrow dynamic viscosity

g \leftarrow gravitational acceleration

k \leftarrow thermal conductivity

λ \leftarrow Stokes hypothesis constant

c_p \leftarrow specific heat, constant pressure

c_v \leftarrow specific heat, constant volume



Vector form

The system (1) can be rewritten in vector form

$$\frac{\partial \mathbf{q}}{\partial t} + \nabla \cdot \mathbf{F}(\mathbf{q}) = \mathbf{S}(\mathbf{q}), \quad (2)$$

for the state variables

$$\mathbf{q} = \begin{pmatrix} \rho \\ \mathbf{u} \equiv \rho \mathbf{u} \\ E \equiv \rho e \end{pmatrix} \leftarrow \begin{array}{l} \text{volume mass density} \\ \leftarrow \text{momentum density} \\ \leftarrow \text{energy density} \end{array} \quad (3)$$



Vector form

The system (1) can be rewritten in vector form

$$\frac{\partial \mathbf{q}}{\partial t} + \nabla \cdot \mathbf{F}(\mathbf{q}) = \mathbf{S}(\mathbf{q}), \quad (2)$$

for the state variables

$$\mathbf{q} = \begin{pmatrix} \rho \\ \mathbf{u} \equiv \rho \mathbf{u} \\ E \equiv \rho e \end{pmatrix} \leftarrow \begin{array}{l} \text{volume mass density} \\ \text{momentum density} \\ \text{energy density} \end{array} \quad (3)$$

where

$$\mathbf{F}(\mathbf{q}) = \begin{pmatrix} \mathbf{u} \\ (\mathbf{u} \otimes \mathbf{u})/\rho + P\mathbf{I}_3 - \boldsymbol{\sigma} \\ (E + P)\mathbf{u}/\rho - (\mathbf{u} \cdot \boldsymbol{\sigma} + k\nabla T) \end{pmatrix},$$

$$\mathbf{S}(\mathbf{q}) = - \begin{pmatrix} 0 \\ \rho g \hat{\mathbf{k}} \\ 0 \end{pmatrix}$$



Application example: Density current

A cold air bubble drops by convection in a neutrally stratified atmosphere.

Its initial condition is defined in terms of the Exner pressure, $\pi(x, t)$, and potential temperature, $\theta(x, t)$, that relate to the state variables via

$$\rho = \frac{P_0}{(c_p - c_v)\theta(x, t)} \pi(x, t)^{\frac{c_v}{c_p - c_v}}, \quad (4a)$$

$$e = c_v\theta(x, t)\pi(x, t) + \mathbf{u} \cdot \mathbf{u}/2 + gz, \quad (4b)$$

where P_0 is the atmospheric pressure.

BCs: no-slip and non-penetration for \mathbf{u} , no-flux for mass and energy densities.



Density current

order: $p = 8$, $P = p + 1$, $q = P + 4$, $\Omega = [0, 6000]^2 \text{ m} \times [0, 3000] \text{ m}$,
elem. resolution: 125 m, Global DOFs: 1 081 665



Conclusions and outlook

- We have demonstrated the use of libCEED with PETSc for the numerical high-order solutions of the full compressible Navier-Stokes equations
- libCEED always welcomes contributors and friendly users
<https://github.com/CEED/libCEED>



Conclusions and outlook

- We have demonstrated the use of libCEED with PETSc for the numerical high-order solutions of the full compressible Navier-Stokes equations
- libCEED always welcomes contributors and friendly users
<https://github.com/CEED/libCEED>



Ongoing and future work:

- Improved non-conforming and mixed-mesh support (integration of PETSc's DMFlex)
- Implicit-explicit (IMEX) and fully implicit time discretizations
- Algorithmic differentiation (AD) of Q-functions
- Ongoing work on CUDA and HIP optimizations



Conclusions and outlook

- We have demonstrated the use of libCEED with PETSc for the numerical high-order solutions of the full compressible Navier-Stokes equations
- libCEED always welcomes contributors and friendly users
<https://github.com/CEED/libCEED>



Ongoing and future work:

- Improved non-conforming and mixed-mesh support (integration of PETSc's DMFlex)
- Implicit-explicit (IMEX) and fully implicit time discretizations
- Algorithmic differentiation (AD) of Q-functions
- Ongoing work on CUDA and HIP optimizations

Acknowledgements: Exascale Computing Project (17-SC-20-SC), NSF

Thank you!



University of Colorado
Boulder