

deal.II on GPU

Final presentation at Eurohack 19, Lugano, Switzerland

Peter Munch¹²³, Momme Allalen⁴, Martin Kronbichler¹²,
Paddy Ó Conbhuí⁴, Prashanth Kanduri⁵

¹deal.II developer

²Institute for Computational Mechanics, Technical University of Munich, Germany

³Institute of Materials Research, Materials Mechanics, Helmholtz-Zentrum Geesthacht, Germany

⁴CFD-Lab, Leibniz Supercomputing Centre, Germany

⁴Irish Centre for High-End Computing, Ireland

⁵CSCS/ ETH Zürich, Switzerland

October 4, 2019

Part 1:

Motivation

Goal:

Optimize the matrix-free deal.II¹ GPU implementation for CEED BP5².

¹General purpose finite element library deal.II: <https://www.dealii.org/>.

²Solve Poisson problem for $p = 5$ and $q = p + 1$ with preconditioned conjugate gradient methods.

Given modules in deal.II:

- ▶ efficient matrix-free implementation
- ▶ GPU support with CUDA-aware MPI-support (with not satisfying performance results)
- ▶ standard preconditioned conjugate gradient method (next slide)

Kronbichler, Kormann, A generic interface for parallel finite element operator application. *Comput. Fluids* 63:135–147, 2012

Kronbichler, Kormann, Fast matrix-free evaluation of discontinuous Galerkin finite element operators. *arXiv:1711.03590*, 2017

M. Kronbichler, K. Ljungkvist. Multigrid for matrix-free high-order finite element computations on graphics processors, *ACM Transactions on Parallel Computing*, 2019

Daniel Arndt and Bruno Turcksin, Tutorial step-64, online available https://www.dealii.org/developer/doxygen/deal.II/step_64.html

Algorithm 1: (Standard) preconditioned CG solver

Data: \mathbf{A} , x_0 , b , \mathbf{M}^{-1}

Result: x

1 [startup];

2 **for** $k = 0, 1, \dots$ **do**

3 $\underline{\mathbf{v}} \leftarrow \underline{\mathbf{A}}\underline{\mathbf{p}}$;

4 $\alpha \leftarrow \gamma / (\underline{\mathbf{p}}^T \underline{\mathbf{v}})$;

5 $\underline{\mathbf{x}} \leftarrow \underline{\mathbf{x}} + \alpha \cdot \underline{\mathbf{p}}$;

6 $\underline{\mathbf{r}} \leftarrow \underline{\mathbf{r}} - \alpha \cdot \underline{\mathbf{v}}$;

7 **if** $\|\underline{\mathbf{r}}\|_2 < \varepsilon$ **then**

8 **return**

9 $\underline{\mathbf{v}} \leftarrow \underline{\mathbf{M}}^{-1} \underline{\mathbf{r}}$;

10 $\beta \leftarrow \gamma$;

11 $\gamma \leftarrow \underline{\mathbf{r}}^T \underline{\mathbf{v}}$;

12 $\beta \leftarrow \gamma / \beta$;

13 $\underline{\mathbf{p}} \leftarrow \underline{\mathbf{v}} + \beta \cdot \underline{\mathbf{p}}$;

- ▶ $\underline{\mathbf{A}}\underline{\mathbf{d}}$ evaluated matrix-free
- ▶ vectors live either on host or device
- ▶ CG algorithm can operate on both types of vectors
- ▶ CPU and GPU implementation of vector operations

- ▶ note: $\underline{\mathbf{M}} := \text{diag}(\underline{\mathbf{A}})$ here!

Algorithm 2: Cache friendly CG preconditioned by a diagonal

Data: \mathbf{A} , x_0 , b , \mathbf{M}^{-1}

Result: x

```

1 [startup];
2 for  $k = 0, 1, \dots$  do
3    $\underline{r} \leftarrow \underline{r} - \alpha \underline{v}$     $\underline{x} \leftarrow \underline{x} + \alpha \underline{p}$     $\underline{p} \leftarrow \underline{\underline{M}}^{-1} \underline{r} + \beta \underline{p}$ ;           // pre
4    $\underline{v} \leftarrow \underline{\mathbf{A}} \underline{p}$ ;                                           // vmult
5    $a \leftarrow \underline{p}^T \underline{v}$    $b \leftarrow \underline{r}^T \underline{r}$    $c \leftarrow \underline{r}^T \underline{v}$    $d \leftarrow \underline{v}^T \underline{v}$    $e \leftarrow \underline{r}^T \underline{\underline{M}}^{-1} \underline{r}$    $f \leftarrow \underline{r}^T \underline{\underline{M}}^{-1} \underline{v}$    $g \leftarrow \underline{v}^T \underline{\underline{M}}^{-1} \underline{v}$ ;           // post
6    $\alpha \leftarrow e/a$ ;
7   if  $\sqrt{b + 2\alpha c + \alpha^2 d} < \epsilon$  then
8      $\underline{x} \leftarrow \underline{x} + \alpha \underline{p}$ ;
9     return
10   $\beta \leftarrow (e - 2\alpha f + \alpha^2 g)/e$ ;

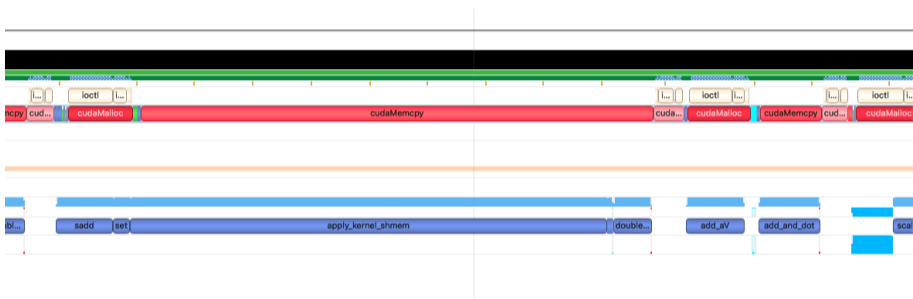
```

Observation: Works very well for CPUs!

Part 2:

Optimization steps

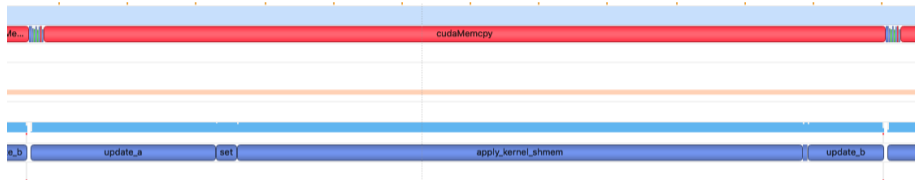
Profile of a single iteration (with NVIDIA Nsight):



Observations:

- + schoolbook-like: easy to understand; working for any preconditioners
- (too) many kernel calls and redundant memory loads

Profile of a single iteration (with NVIDIA Nsight):



Observations:

- + 3 kernel calls³
- limited to diagonal preconditioner

³We eliminated "set" and merged into "update_a" in a next step (not shown here) with reduced memory access.

Profile of a single iteration (with NVIDIA Nsight):



Observations:

- No overlap of communication and computation.

Overlapping communication and computation (multi node/single GPU)



Profile of a single iteration (with NVIDIA Nsight):



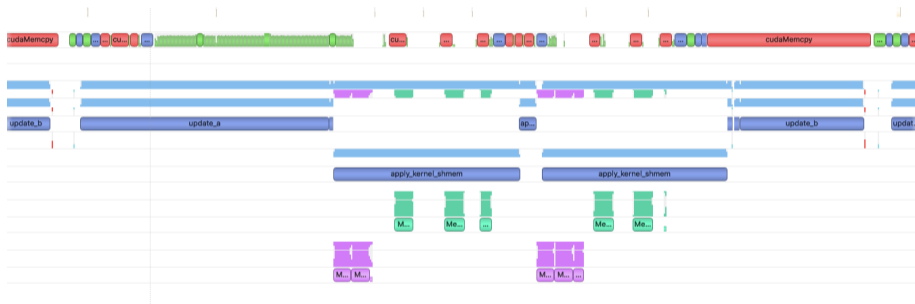
Split up matrix-vector multiplication into three parts → Observations:

- + Overlap of (actual) communication and computation.
- MPI (MPI_Isend, MPI_Waitall) memcopies data (DtoH and HtoD), which blocks communication.

Overlapping MPI memcopy with computation (multi node/single GPU)



Profile of a single iteration (with NVIDIA Nsight):



Move first and last part of matrix-vector multiplication on to a new streams → Observations:

- + Implicit synchronization (by default stream).
- + Perfect overlap.

- ▶ work in progress
- ▶ hope: better utilization of the GPU and the network card

The final code can be found online:

- ▶ <https://github.com/peterrum/deal-and-ceed-on-gpu>
- ▶ <https://github.com/peterrum/dealii/tree/dealii-on-gpu>

Part 3:
Results

Achieved throughput:

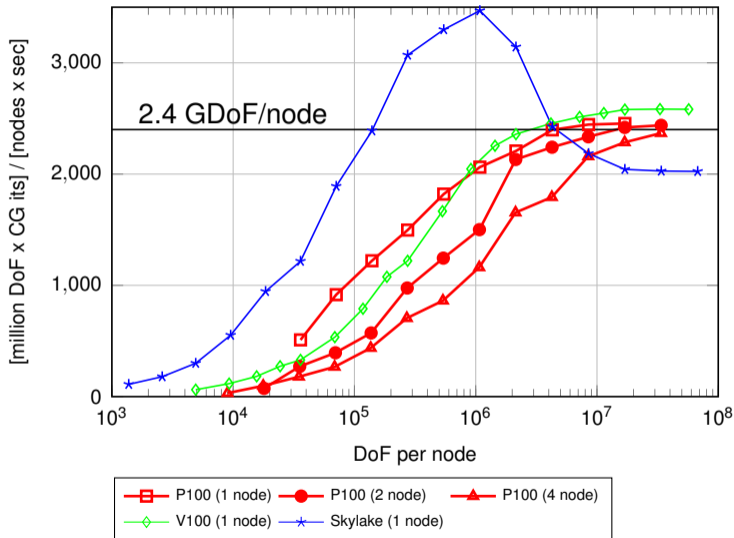
- ▶ **Pascal (4×): 7.6 GDoF**

Total peak performance:

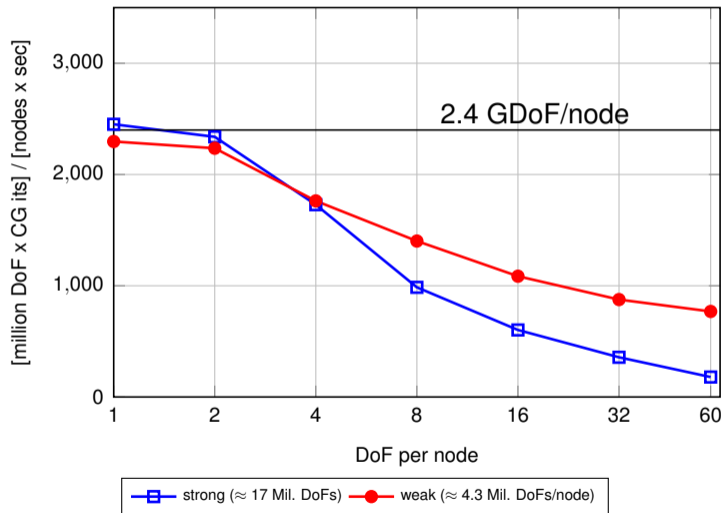
- ▶ **Pascal: 4.7 TFlop/s**
- ▶ Volta: 7.8 TFlop/s
- ▶ Skylake: 4.2 TFlop/s

Maximum memory bw:

- ▶ **Pascal: 732 GB/s**
- ▶ Volta: 900 GB/s
- ▶ Skylake: 191 GB/s



Note: Skylake/Pascal run merged PCG. Volta runs standard implementations!



Variant 1:

- ▶ do not merge coefficient

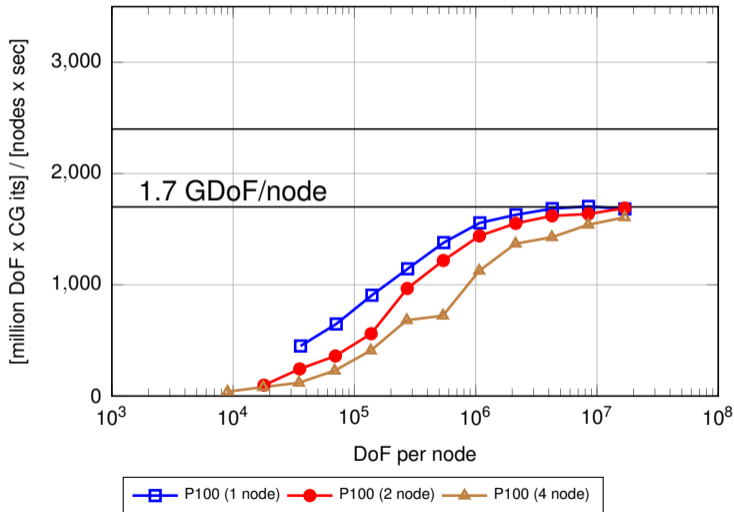
$$(J_q^{-1} |J_q| w_q J_q^{-T})$$

in quadrature points

- ▶ instead: load J_q , $|J_q|$
- ▶ i.e.: load 10 vs. 6 doubles
- ▶ (additional work)

Observation:

- ▶ **Performance drop higher than expected!**

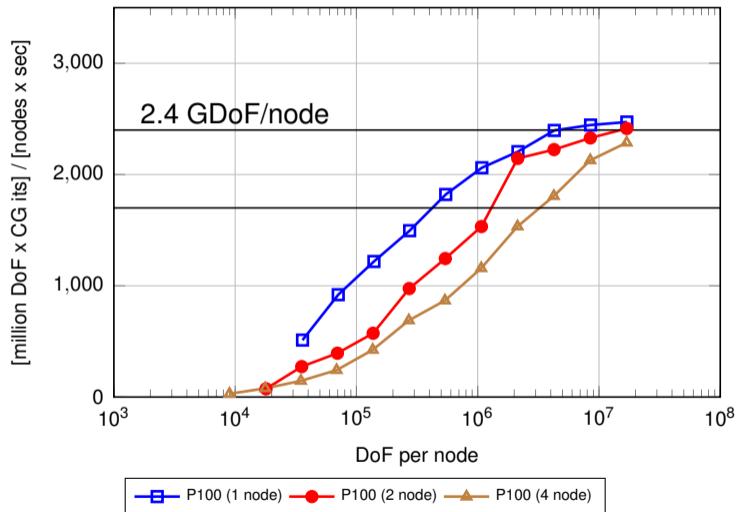


Variant 2:

- ▶ integrate in Gauss quadrature points
- ▶ requires interpolation
- ▶ additional 2 · dim sum-factorization sweeps

Observation:

- ▶ **Basis transformation comes for free!**



Part 4:

Conclusions and lessons learned

- ▶ GPU programming is fun (although we did not have to do much), except the redundant conversion of block- and thread-id
- ▶ similar optimization strategies work for CPU and GPU (e.g. minimize data access overlapping communication and computation)
- ▶ given a heavily optimized FEM code (matrix-free exploiting caches and SIMD), the usage of GPUs will not lead to a significant speed-up

- ▶ Step-by-step integration of the new features in deal.II.
- ▶ Extend the functionality of MatrixFree (GPU) so that it matches MatrixFree (CPU).
 - ▶ geometric multigrid and hybrid multigrid methods⁴
 - ▶ discontinuous Galerkin methods
 - ▶ over-integration $p+1 \ll q$
- ▶ Investigation of compute nodes with multiple GPUs.
- ▶ Build a minimal build configuration for GPU-only runs.

⁴ P. Munch, An efficient hybrid multigrid solver for high-order discontinuous Galerkin methods, Master's thesis, 2018
(online available: <https://mediatum.ub.tum.de/node?id=1514962>)

Thanks to the organizers and the mentors!