



WWU
MÜNSTER



FRIEDRICH-ALEXANDER
UNIVERSITÄT
ERLANGEN-NÜRNBERG
TECHNISCHE FAKULTÄT



UNIVERSITÄT
HEIDELBERG
ZUKUNFT
SEIT 1386

Recent advances for code generation in the *HPC²SE* project

Hardware- and Performance-aware Codegeneration for
Computational Science and Engineering

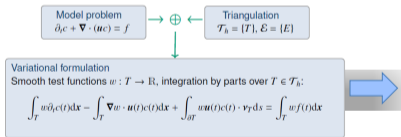
Engwer, Fahlke, Koch, Bastian, Kempf, Heß, Gorlatch, Fey, Rüde, Köstler, Hönig

18. October 2019



Goals for Code Generation

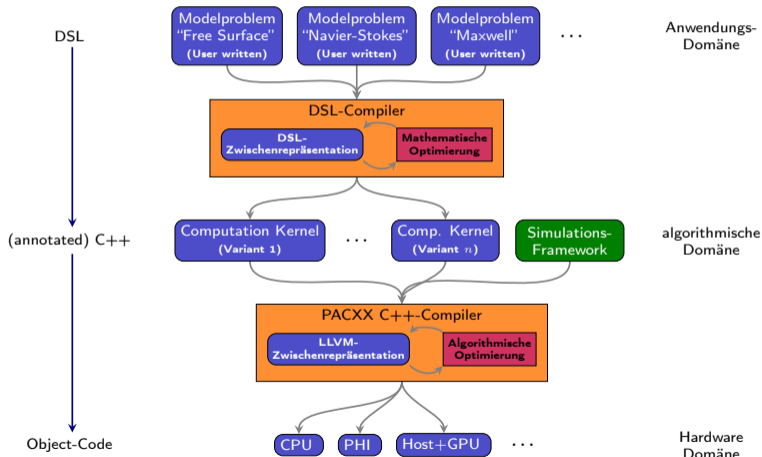
- ▶ Fast problem transformations (productivity)
 - ➔ Math to Code
- ▶ Fast implementations (performance)
 - ➔ Efficient use of hardware resources
- ▶ Fast code transformations (portability)
 - ➔ For algorithms and different platforms



```
#include "MultiGrid/MultiGrid.h"
void Smoother_4O {
  exhsolutionData_4O();
  #pragma omp parallel for schedule(static) num_threads(8)
  for (int fragmentIdx = 0; fragmentIdx < 8; ++fragmentIdx) {
    if (isValidForSubdomain[fragmentIdx][0]) {
      for (int y = iterationOffsetBegin[fragmentIdx][0];
           y < (iterationOffsetEnd[fragmentIdx][0] + 17); y += 1)
        for (int x = iterationOffsetBegin[fragmentIdx][0];
             x < (iterationOffsetEnd[fragmentIdx][0] + 17); x += 1)
          slottedFieldData_Solution[0][fragmentIdx][0][((y*10)+x)+(x+1)] =
            ((slottedFieldData_Solution[0][fragmentIdx][0][((y*10)+x)+(x+1)] +
              (((1.0e+00)/fieldData_LaplCoeff[fragmentIdx][0][((y*17)+x)])*8.0e-01)
              *(fieldData_RHS[fragmentIdx][0][((y*17)+x)]
                - (((fieldData_LaplCoeff[fragmentIdx][0][((y*17)+x)]
                  *(slottedFieldData_Solution[0][fragmentIdx][0][((y*10)+10)+(x+1)]))
                  + (fieldData_LaplCoeff[fragmentIdx][0][((y*7)+28)+(x+1)]))
                  + (fieldData_LaplCoeff[fragmentIdx][0][((y*10)+10)+(x+1)]))
                  + (fieldData_LaplCoeff[fragmentIdx][0][((y*17)+867)+x])
                  + (fieldData_Solution[0][fragmentIdx][0][((y*10)+38)+(x+1)]))
                  + (fieldData_LaplCoeff[fragmentIdx][0][((y*17)+1156)+x])
                  *(slottedFieldData_Solution[0][fragmentIdx][0][((y*10)+(x+1))))))));
    }
  }
}
```



HPC²SE code generation approach

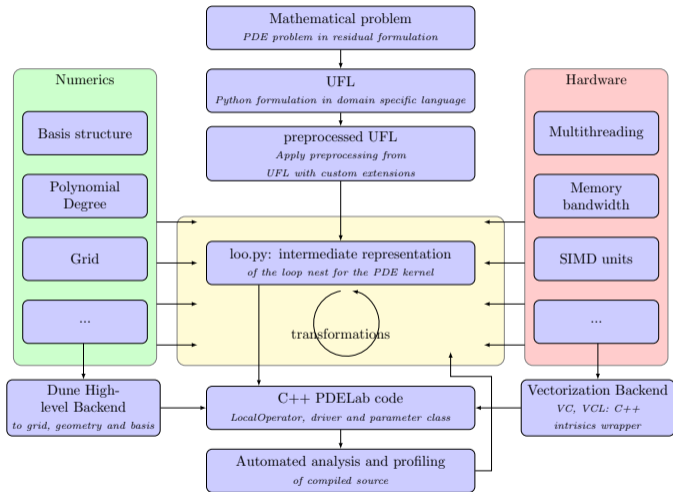


Outline

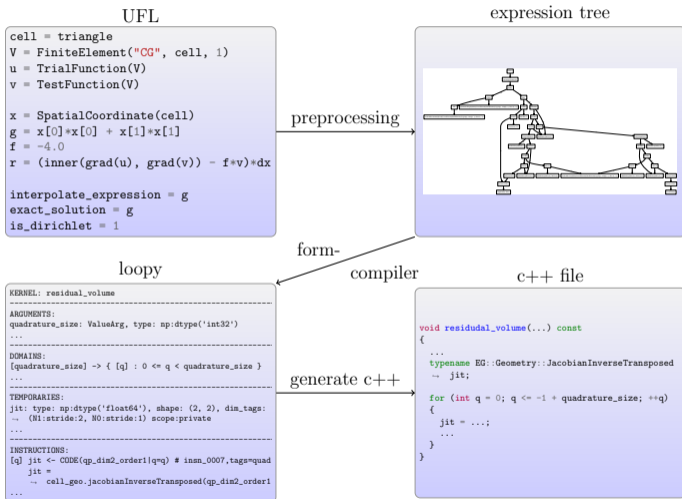
- 1 Our Code Generation Approach
- 2 Higher-order DG
- 3 Block-structured meshes
- 4 Stencil code generation
- 5 Platform portability?!
- 6 Discussion

1 Our Code Generation Approach

Workflow



Workflow



Current State

Low-level optimization

- ▶ Performance optimizations on intermediate representation (eg. loop reordering in tensor contractions)
- ▶ Built in autotuning approach using micro benchmarks for finding good performance transformations
- ▶ Vectorization strategies for good SIMD parallelism

Code generation options

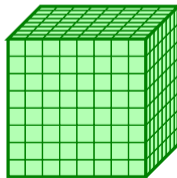
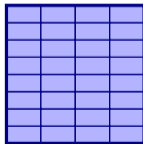
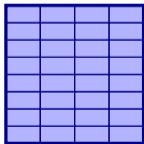
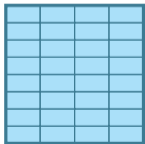
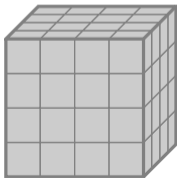
- ▶ Sumfactorization for High order DG
 - ➔ reduces algorithmic complexity
- ▶ Block-structured meshes for low order CG
 - ➔ Increase of arithmetic intensity by handling multiple elements in local operator
- ▶ Stencil kernels on structured meshes
 - ➔ On-the-fly smoothers for MG (➔ EXA-Stencil)

2 Higher-order DG

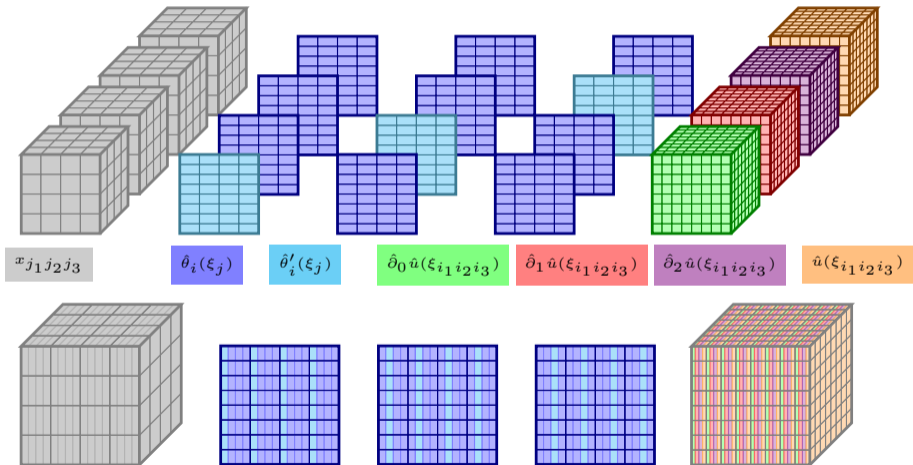
Evaluation of $\hat{\delta}_0 \hat{u}$ in 3D using sum factorization

$$\hat{\delta}_0 \hat{u}(\xi_{i_1 i_2 i_3}) = \underbrace{\sum_{j_3 \in J^{(3)}} A_{i_3, j_3}^{(3)}}_{y_{i_3 i_2 i_1}} \underbrace{\sum_{j_2 \in J^{(2)}} A_{i_2, j_2}^{(2)}}_{y_{j_2 j_3 i_1}} \underbrace{\sum_{j_1 \in J^{(1)}} A_{i_1, j_1}^{(1)} x_{j_1 j_2 j_3}}_{y_{j_2 j_3 i_1}}$$

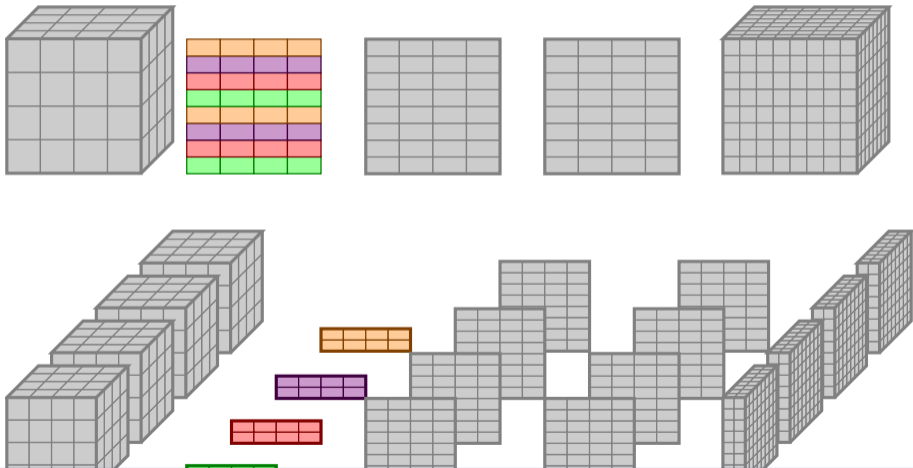
$A_{i,j}^{(1)} = \hat{\theta}'_i(\xi_j^{(1)})$: Derivatives of the 1D Basis at the 1D quadrature points. $A_{i,j}^{(k)} = \hat{\theta}_i(\xi_j^{(k)})$, $k > 1$: Evaluations of the 1D Basis at the 1D quadrature points.



Vectorization by Loop Fusion

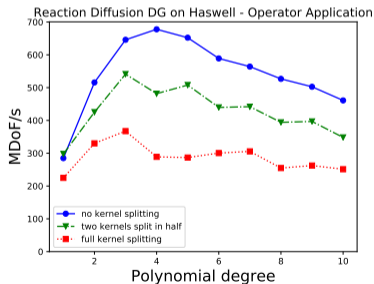
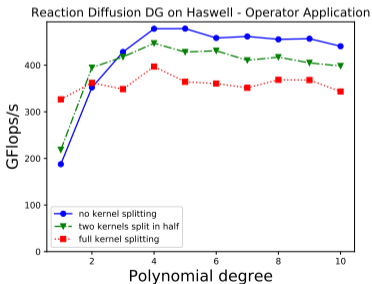


Parallelism through Loop Splitting



Performance on Intel Haswell

- ▶ Intel Xeon processor E5-2698 v3
- ▶ 2×16 cores, 2.3 GHz
- ▶ Theoretical node peak performance: 1.17 TFlops/s
- ▶ Operator application of DG-Method for a diffusion reaction problem with full diffusion tensor.



Peak performance?

Haswell

- ▶ Around 50% of peak performance for polynomial degree $k > 2$
- ▶ Volume integrals around 70% of peak performance
- ▶ Skeletons around 40% of peak performance (lower flop per byte ratio)

Skylake

- ▶ Around 30% of peak performance for polynomial degree $k > 2$
- ▶ Around 70% of peak for volume integrals but room for improvement on skeletons
- ▶ Up to 850 GFlops/s compared to 500 GFlops/s for Haswell

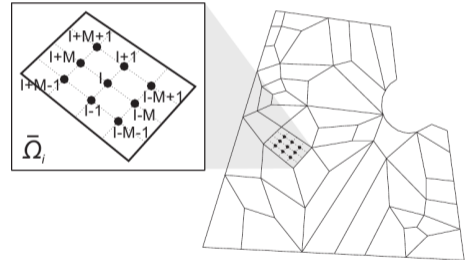
3 Block-structured meshes

Locally block-structured FEM

Goal: Generate efficient lower order FEM code

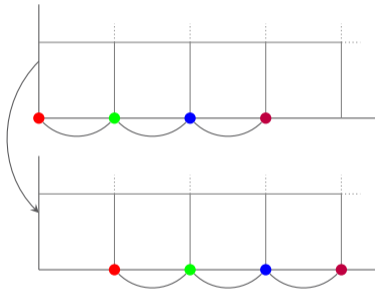
Idea:

- ▶ Coarse grid of macro elements
 - ▶ Refine macro-Elements in k Micro-Elements
 - ▶ Substructure managed within the local kernel
- increased arithmetic intensity



Vectorization

- ▶ Vectorize over neighboring micro-elements
- ▶ Automatically generated as transformation within the Loo.py IR
- ▶ Example vector width 4:
→ local kernel speedup of $\sim 3.5x$



Performance - Matrix free operator application

Linear elasticity (LE)

Unstructured 3D mesh

Local kernel with high arithmetic intensity

40MB DoFs per kernel

Poisson (P)

Axiparallel strukturiertes 2D mesh

Local kernel with low arithmetic intensity

95MB DoFs per kernel

Problem	Variant	MDof/s	Speedup vs $k = 1$
LE	$k = 8$	2.1	2.2x
	$k = 8 + \text{SIMD4}$	7.1	7.8x
P	$k = 16$	18.4	5.6x
	$k = 16 + \text{SIMD4}$	41.3	12.6x

Performance - Matrix free operator application

Linear elasticity (LE)

Unstructured 3D mesh

Local kernel with high arithmetic intensity

40MB DoFs per kernel

Low data transfer overhead

→ low speedup due to block-structuring

High arithmetic intensity

→ high speedup due to vektorization

Poisson (P)

Axiparallel strukturiertes 2D mesh

Local kernel with low arithmetic intensity

95MB DoFs per kernel

High data transfer overhead

→ hoher speedup due to block-structuring

Low arithmetic intensity

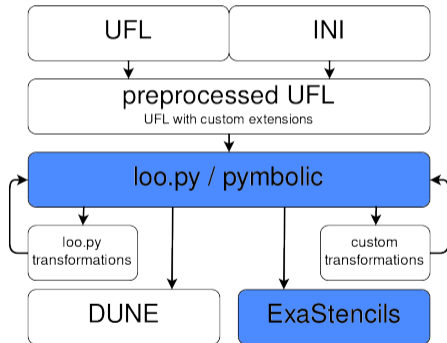
→ low speedup due to vektorization

4 Stencil code generation

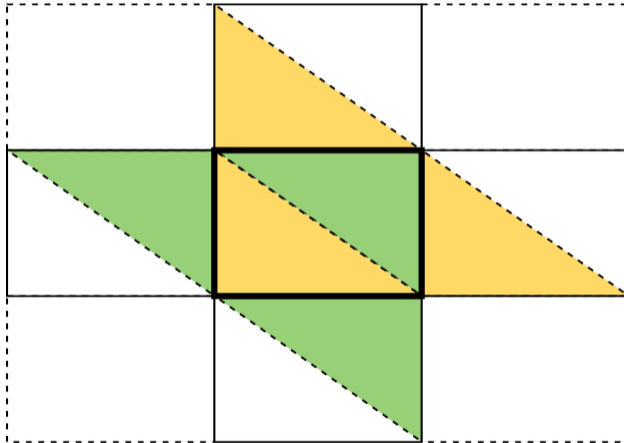
Stencil generation for HPC platforms

ExaStencils

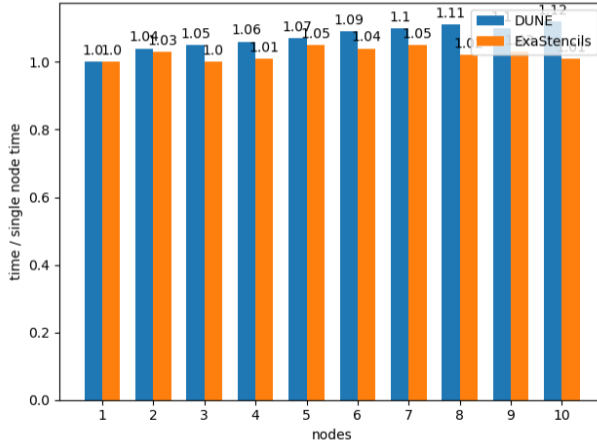
- ▶ Whole program generation framework
- ▶ Focus on Multigrid methods and stencil codes
- ▶ Own DSL: ExaSlang
- ▶ Target Code: C++/CUDA
- ▶ Automatically introduces SIMD, OpenMP, MPI
- ▶ Heterogeneous computation



DG Stencil for ExaStencils



Scaling on SuperMUC-NG



5 Platform portability?!

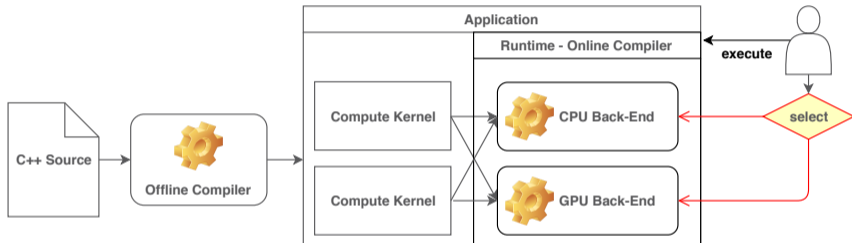
PACXX: Platform portable C++ compiler

- ▶ PACXX is used as a complete compiler
- ▶ Drop-in replacement for Clang
- ▶ Compilation result contains »Runtime« + all support routines

Objectives

- ▶ Unified programming for parallel architectures in C++.
- ▶ Platform independence by modular design.
- ▶ “Black-box” parallelisation.

Unified programming for parallel architectures in C++



- ▶ Inspired by CUDA and OpenCL.
- ▶ »Kernel« / »Devices« paradigm.
- ▶ One codebase for »Host« and »Device«.
- ▶ Based on C++14 Lambdas.

Example: MatMul

```
auto& exec = Executor::get(0);

auto& dev_a = exec.allocate<double>(matrix_size);
auto& dev_b = exec.allocate<double>(matrix_size);
auto& dev_c = exec.allocate<double>(matrix_size);

dev_a.upload(a, matrix_size);
dev_b.upload(b, matrix_size);
dev_c.upload(c, matrix_size);

auto pa = dev_a.get();
auto pb = dev_b.get();
auto pc = dev_c.get();

auto matMultKernel = [=](auto &config)
{
    auto column = config.get_global(0);
    auto row    = config.get_global(1);
    double val = 0;
    for (unsigned i = 0; i < width; ++i)
        val += pa[row * width + i] * pb[i * width + column];
    pc[row * width + column] = val;
};

exec.launch(matMultKernel, {{width/threads , width} , {threads , 1}});

dev_c.download(c, matrix_size);
```

Example: MatMul

```
auto& exec = Executor::get(0);

auto& dev_a = exec.allocate<double>(matrix_size);
auto& dev_b = exec.allocate<double>(matrix_size);
auto& dev_c = exec.allocate<double>(matrix_size);

dev_a.upload(a, matrix_size);
dev_b.upload(b, matrix_size);
dev_c.upload(c, matrix_size);

auto pa = dev_a.get();
auto pb = dev_b.get();
auto pc = dev_c.get();

auto matMultKernel = [=](auto &config)
{
    auto column = config.get_global(0);
    auto row    = config.get_global(1);
    double val = 0;
    for (unsigned i = 0; i < width; ++i)
        val += pa[row * width + i] * pb[i * width + column];
    pc[row * width + column] = val;
};

exec.launch(matMultKernel, {{width/threads , width} , {threads , 1}});

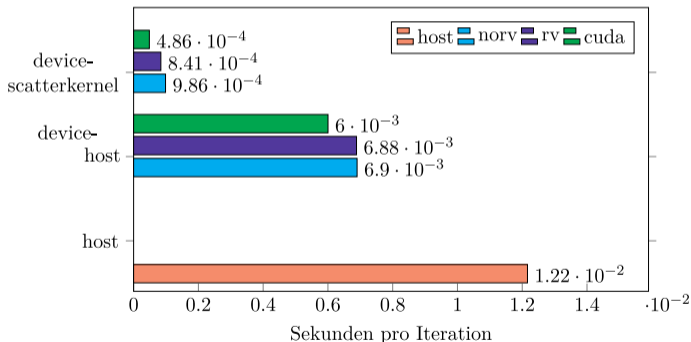
dev_c.download(c, matrix_size);
```

Using PACXX GPU backend for grid-based methods

Intel Xeon CPU E5-1620 v2 vs. Nvidia Tesla K20c

refinement=6 (24833 cells)

- ▶ Models: »Linear Elasticity« and »Poisson«
- ▶ Reimplement certain DUNE components “GPU friendly”
- ▶ Reformulate operator/assembler similar to Map/Reduce
- ▶ Split into two separate kernels





WWU
MÜNSTER



Advanced Solvers for modern Architectures 2019

7th Applied Mathematics Symposium Münster
November 11–13, 2019 | Münster, Germany

uni-muenster.de/AMM/num/solvers2019/

Speakers

Andrew T. Barker (LLNL)

Erin Carson (Charles University)

Andreas Frommer (University of Wuppertal)

Laura Grigori (INRIA)

Eike Müller (University of Bath)

Nicole Spillane (École Polytechnique)

Wim Vanroose (University of Antwerp)

Stefano Zampini (KAUST)

Discussion

- ▶ Flexible code generation pipeline
- ▶ Mathematically increased arithmetic intensity via local structure
 - ▶ Higher-order DG
 - ▶ Block-structured refinement
- ▶ Generating preconditioners
- ▶ Platform portability
 - ▶ GPU friendly mesh data structures
 - ▶ Map-reduce style assembly
 - ▶ PACXX concepts not generic enough for our optimizations
- ▶ Open Questions:
 - ▶ Transfer to real application
 - ▶ Finalize work on DD preconditioners
 - ▶ Fully integrate with PACXX

Discussion

- ▶ Flexible code generation pipeline
- ▶ Mathematically increased arithmetic intensity via local structure
- ▶ Generating preconditioners
 - ▶ Stencil generation
 - ▶ *Non-verlapping DD methods (WIP)*
- ▶ Platform portability
 - ▶ GPU friendly mesh data structures
 - ▶ Map-reduce style assembly
 - ▶ PACXX concepts not generic enough for our optimizations
- ▶ Open Questions:
 - ▶ Transfer to real application
 - ▶ Finalize work on DD preconditioners
 - ▶ Fully integrate with PACXX

Questions?