PyFR and GiMMiK on Intel KNL

¹Department of Aeronautics & Astronautics, Stanford University ²Departments of Aeronautics and Computing, Imperial College London ³Intel Corporation

F.D. Witherden¹, J.S. Park², A. Heinecke³, P. Kelly², P.E. Vincent² and A. Jameson¹





• Interested in simulating unsteady, turbulent, flows.





Motivation: Turbulent Flows







unstructured grids with explicit time stepping.

The PyFR Framework

DPyFR

• Uses high-order flux reconstruction (FR) to solve the compressible Navier-Stokes equations on mixed

The PyFR Framework

• Performance portable across a range of platforms.

• Finalist for the 2016 Gordon Bell Prize.



• Existing support for KNC based around offloading via pyMIC.

The PyFR Framework

• Python outer layer.

Python Outer Layer (Hardware Independent)

- Setup
- Distributed memory parallelism
- Outer 'for' loop and calls to hardware specific kernels

• Need to generate hardware specific kernels.

Python Outer Layer (Hardware Independent)

- Setup
- Distributed memory parallelism
- Outer 'for' loop and calls to hardware specific kernels

• In FR two types of kernel are required.

Python Outer Layer (Hardware Independent)

- Setup
- Distributed memory parallelism
- Outer 'for' loop and calls to hardware specific kernels

PyFR

Matrix Multiply Kernels

• Data interpolation/ extrapolation etc.

Point-Wise Nonlinear Kernels

• Flux functions, Riemann solvers etc.

• Matrix multiplications are quite simple.

Python Outer Layer (Hardware Independent)

- Setup
- Distributed memory parallelism
- Outer 'for' loop and calls to hardware specific kernels

PyFR

Matrix Multiply Kernels

• Data interpolation/ extrapolation etc.

Point-Wise Nonlinear Kernels

• Flux functions, Riemann solvers etc.

Call GEMM

• For the point-wise nonlinear kernels we use a DSL.

Python Outer Layer (Hardware Independent)

- Setup
- Distributed memory parallelism
- Outer 'for' loop and calls to hardware specific kernels



• Kernels are generated and compiled at start-up.

Python Outer Layer



• Which may then be called by the outer layer.



Matrix Multiplications in PyFR

• Multiplications are of the **block-by-panel** variety:



• where $N \sim 10^5$ with $N \gg (M, K)$ and **A is constant**.



GEMM in PyFR

• On x86 S/DGEMM has three kernels providers.



- Flow over a cylinder at Re = 3,900 and Ma = 0.2.



Initial Results

• Quadratically curved hexahedral mesh with $N_E = 118,820$.





Initial Results

• PyFR 1.4.0: K40c (cuBLAS) vs KNL 7250F (MKL).



Time per DOF per RK stage / ns

Initial Results

- - Surprising!
- vectorisation.

Initial Results

• Profiling indicates **point-wise kernels** are the bottleneck.

• Must therefore rethink our data layout and push for further

Data Layout 101

- Three main layouts:
 - AoS
 - SoA (used by PyFR 1.4.0)
 - AoSoA(k)



struct float rho; float rhou; float E; } data[NELES];

AoS



Memory

- Cache and TLB friendly.
- Difficult to vectorise.





struct float E[NELES]; data;



float rho[NELES]; float rhou[NELES];





- Trivial to vectorise.
- Can put pressure on TLB and/or hardware pre-fetchers.



AoSoA(k = 2)

struct
{
 float rho[k];
 float rhou[k];
 float E[k];
 data[NELES / k];



Memory

- Can be vectorised efficiently for suitable k.
- Cache and TLB friendly.

The Goldilocks solution

• ... albeit at the cost of messy indexing

• ...also requires coaxing for compilers to vectorise.

AoSoA(k = 2)

AoSoA(k = 2)

- Here PyFR's DSL pays dividends.
- Iteration and indexing are hidden from kernels.
- Can therefore change data structures with ease.

AoSoA(k) Results • PyFR 1.4.0 KNL (MKL) vs PyFR 1.5.0 KNL (MKL). p = 1p = 2 p = 3p = 4



12 10 8 6 Time per DOF per RK stage / ns

AoSoA(k) Results

• Comparing with PyFR 1.4.0 K40c (cuBLAS)



Time per DOF per RK stage / ns

AoSoA(k) Results

- Large performance boost.
 - KNL now outperforms a K40c in the dense regime.
 - Other backends improve, too, but only slightly.

Sparsity Exploitation

- Hexahedral operator matrices are sparse.
- Can therefore use GiMMiK/libxsmm.



Sparsity Exploitation

- All with PyFR 1.5.0
 - K40c (GiMMiK)
 - KNL (GiMMiK)
 - KNL (libxsmm)



Summary

- Performance is promising.
- libxsmm.

• Lots of room for improvement in GiMMiK and

Acknowledgements

• Air Force Office of Scientific Research for support under grant FA9550-14-1-0186.





STANFORD UNIVERSITY

Backup Slides

- Normally start up time is a non-issue for PyFR.
- days.

• A few minutes for a simulation which will run for hours to

• Time is split roughly equally between (i) running serial Python code;

(ii) using **ICC** to compile run-time generated kernels.

- Difficult to improve
 - ... Python is virtually impossible to JIT
 - ...and due to the GIL does not benefit from multi-threading.

- But did add a kernel cache.
- Huge reduction in time for p = 4 cylinder case with GiMMiK.

