

EVOLUTIONARY RE-ENGINEERING OF MULTI-PHYSICS INDUSTRIAL HPC APPLICATIONS WITH OP-DSLs

Gihan Mudalige

Royal Society Industry Fellow

Associate Professor (Reader) in High Performance Computing

Department of Computer Science, University of Warwick

g.mudalige@warwick.ac.uk

Joint work with:

Istvan Reguly @ PPCU,

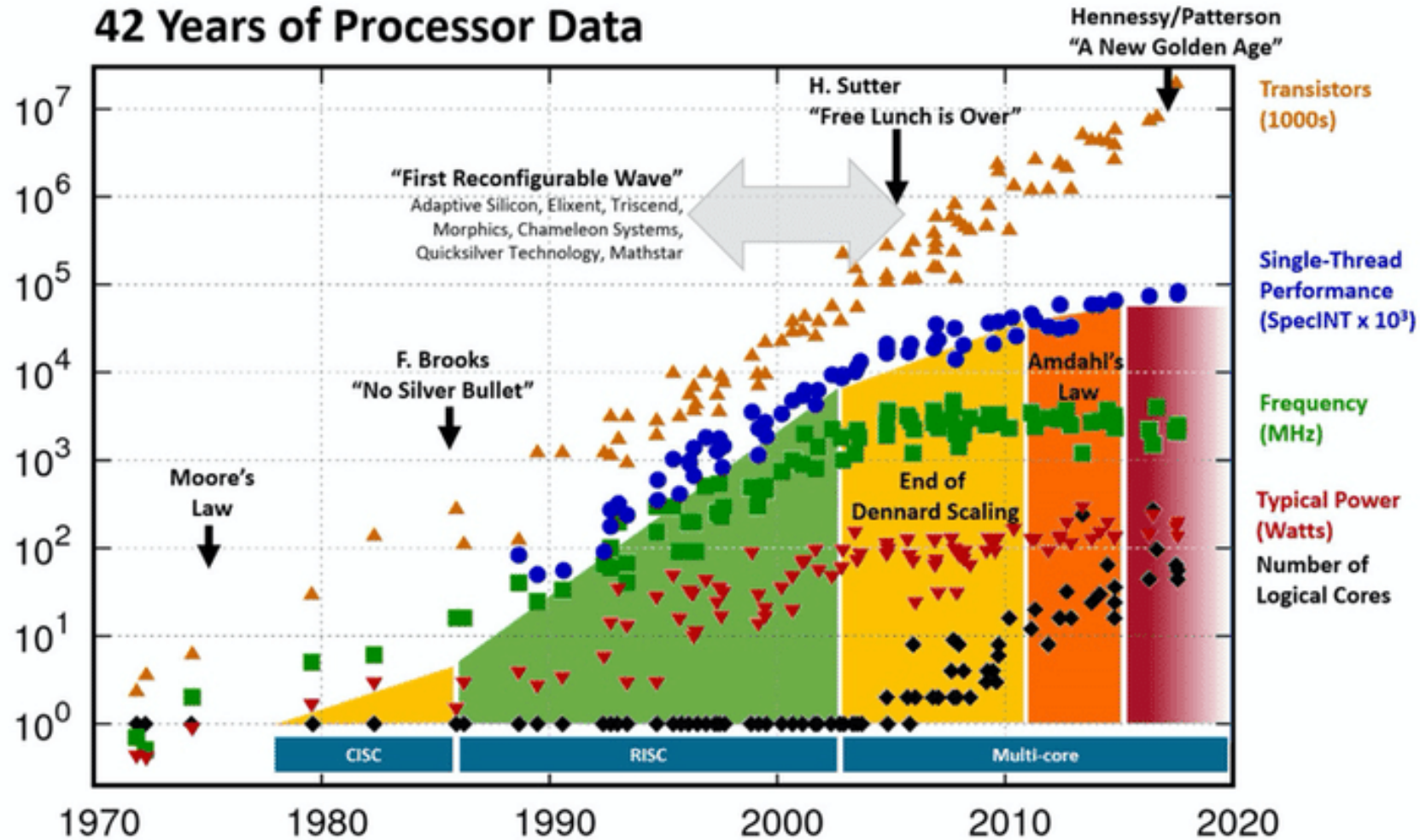
Kamalavasan Kamalakannan, Arun Prabhakar, Archie Powell and others at the HPSC group @ Warwick

Neil Sandham and team @ Southampton, Dario Amirante @ Surrey,

Mike Giles @ Oxford, Paul Kelly and many more @ Imperial College London,

Rolls-Royce plc., NAG, UCL, STFC, IBM and many more industrial and academic collaborators.

SINGLE THREAD SPEEDUP IS DEAD – MUST EXPLOIT PARALLELISM



Hennessy and Patterson, Turing Lecture 2018, overlaid over "42 Years of Processors Data"

<https://www.karlsruhp.net/2018/02/42-years-of-microprocessor-trend-data/>; "First Wave" added by Les Wilson, Frank Schirrmeyer

Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten

New plot and data collected for 2010-2017 by K. Rupp

THE HAIL MARY PASS !



<http://www.themike.com/mikes-free-football-comic-book-hail-mary-pass/>

“The semiconductor industry threw the equivalent of a Hail Mary pass when it switched from making microprocessors run faster to putting more of them on a chip - doing so without any clear notion of how such devices would in general be programmed.”

David Patterson, University of California - Berkeley 2010

DIVERSE HARDWARE LANDSCAPE – COMPOUNDED BY THE RACE TO EXASCALE !

❑ Traditional CPUs

- Intel, AMD, ARM, IBM
- multi-core (> 20 currently)
- Deep memory hierarchy (cache levels and RAM)
- longer vector units (e.g. AVX-512)

❑ GPUs

- NVIDIA (A100), AMD (MI200) , Intel (Xe GPUs)
- Many-core (> 1024 simpler SIMT cores)
- CUDA cores, Tensor cores
- Cache, Shared memory, HBM (3D stacked DRAM)

❑ Heterogeneous Processors

- Different core architectures over the past few years
- ARM big.LITTLE
- NVIDIA Grace.Hopper

❑ XeonPhi (discontinued)

- Many-core – based on simpler x86 cores
- MCDRAM (3D stacked DRAM)

❑ FPGAs

- Dominated by Xilinx and Intel
- Various configurations
- Low-level language / HLS tools for programming
- Significant energy savings

❑ DSP Processors

- Phytium / The Chinese Matrix2000 GPDSP accelerator (Yet to be announced Chinese Exascale system ?)

❑ TPUs, IPU's

❑ Quantum ?

BUT .. EVEN MORE DIVERSE WAYS TO PROGRAMMING THEM !

OpenMP,
SIMD,
CUDA, OpenCL,
OpenMP4.0, OpenACC,
SYCL/OneAPI,
HIP/ROCM,
MPI, PGAS
Task-based (e.g Legion)
and others

- ❑ Open standards (e.g OpenMP, SYCL) – so far have not been agile to catch up with changing architectures
- ❑ Proprietary models (CUDA, OpenACC, ROCm, OneAPI) – restricted to narrow vendor specific hardware
- ❑ Need different code-paths/parallelization schemes to get the best performance
 - ❑ E.g. Coloring vs atomics vs SIMD vs MPI vs Cache-blocking tiling for unstructured mesh class of applications
- ❑ What about legacy codes ? There is a lot of FORTRAN code out there !

SOFTWARE CHALLENGE – A MOVING TARGET

- ❑ What would an Exa-scale machine architecturally look like ?

- ❑ Each new platform requires new performance tuning effort
 - Deeper memory/cache hierarchies and/or shared-memory (non-coherent)
 - Multiple (heterogeneous) memory spaces (device memory/host memory)
 - Complex programming skills set needed to extract best performance on the newest architectures

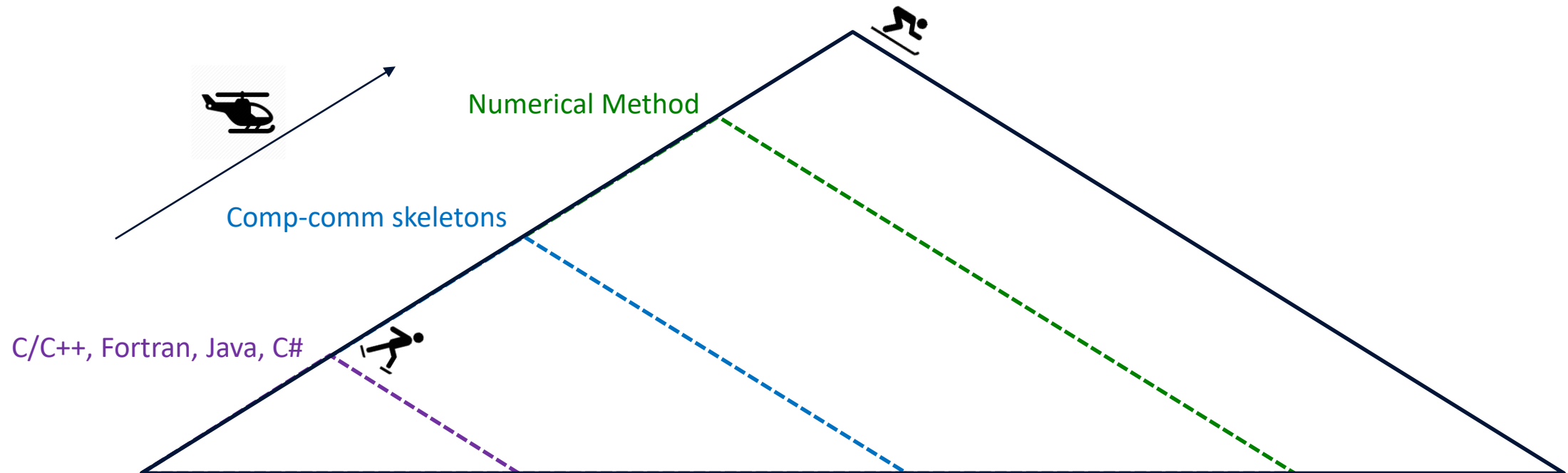
- ❑ Not clear which architectural approach is likely to *win* in the long-term
 - Cannot be re-coding applications for each *new* type of architecture or parallel system
 - Nearly impossible for re-writing legacy codes

- ❑ Need to future-proof applications for their continued performance and portability
 - If not – significant loss of investment : applications will not be able to make use of emerging architectures

OUTLINE

- ☐ Motivation ✓
- ☐ OP-DSLs
- ☐ Evolving Hydra to OP2 Hydra
- ☐ Challenges and Lessons Learnt
- ☐ Conclusions

DOMAIN SPECIFIC ABSTRACTIONS



Adapted from: *Synthesis versus Analysis: What Do We Actually Gain from Domain-Specificity?* Keynote talk at the LCPC 2015. Paul H. J. Kelly (Imperial College London)

- ❑ Rise the abstraction to a **specific domain** of variability
- ❑ Concentrate on a narrower range (class) of computations
 - Computation-Communications skeletons - Structured-mesh, Unstructured-mesh, ... 7 Dwarfs [Colella 2004] ?
 - (higher) Numerical Method - PDEs, FFTs, Monte Carlo ...
 - (even higher) Specify application requirements, leaving implementation to select radically different solution approaches

- ❑ Separation of Concerns (... back in 2010 !)
 - Specify the problem – not the implementation
 - Leverage the best implementation for the target context
 - Can be many contexts - hardware, programming model, parameters etc.

- ❑ Domain Specific API
 - Get application scientists to pose the solution using domain specific constructs – provided by the API
 - Handling data done *only* using API – contract with the user

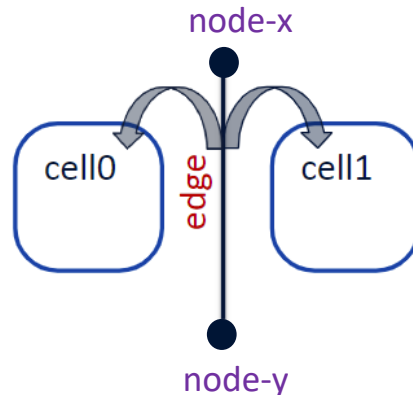
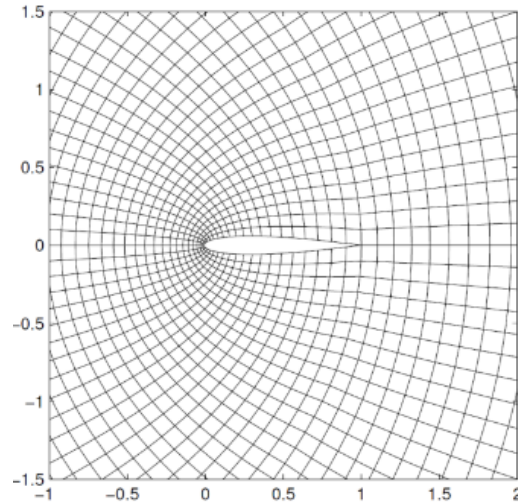
- ❑ Restrict writing code that is difficult (for the compiler) to reason about and optimize
 - “*OP2 and OPS are a straightjacket*” – Mike Giles
 - Build in safe guards so that user cannot write bad code !

- ❑ Implementation of the API left to a lower level
 - Target implementation to hardware – can use best optimizations
 - Automatically generate implementation from specification for the context
 - Exploit domain knowledge for better optimisations - reuse what we know is best for each context



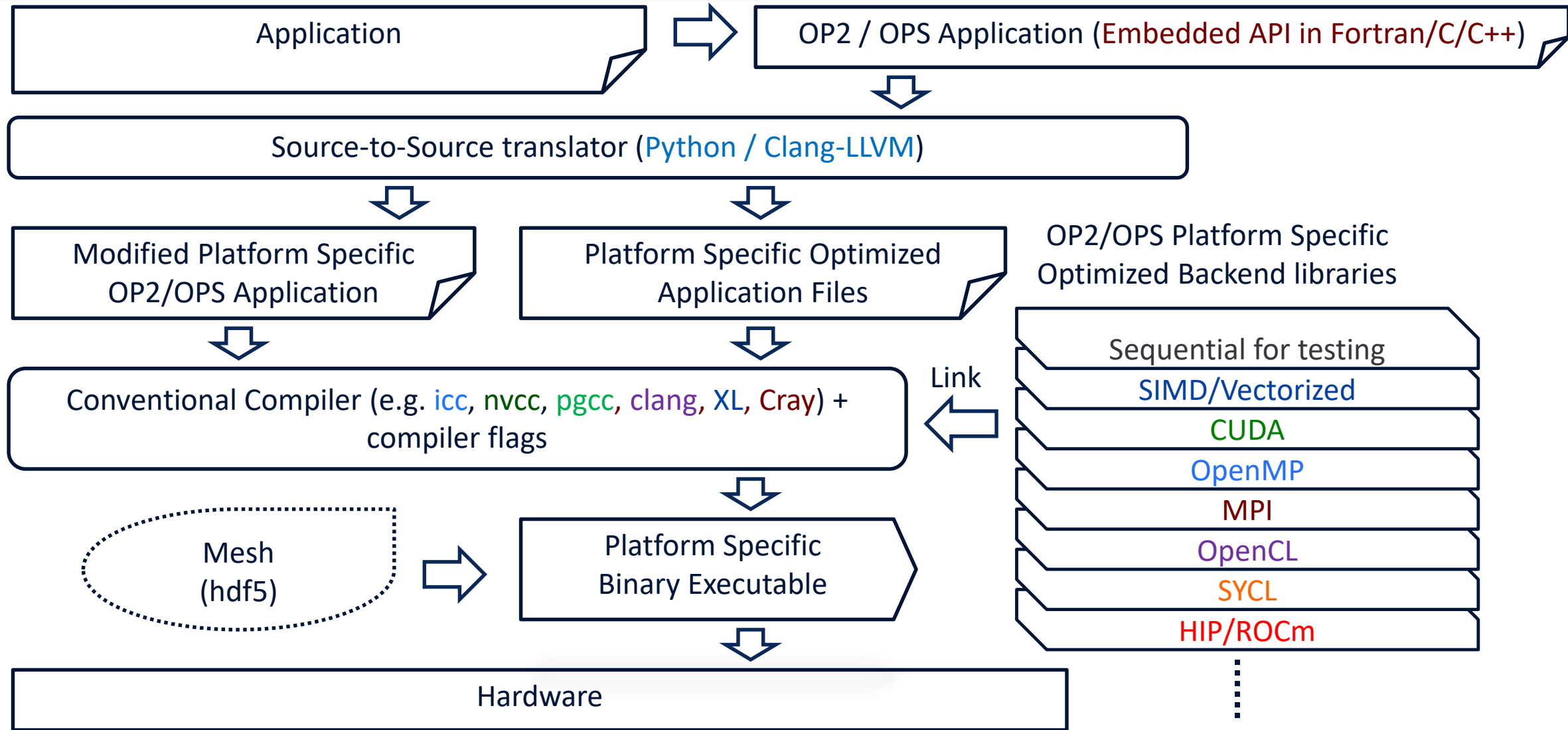
OP2 API - EXAMPLE

```
1 ! Declaring the mesh with OP2
2 ! sets
3 call op_decl_set(nnode,nodes,'nodes')
4 call op_decl_set(nedge,edges,'edges')
5 call op_decl_set(ncell,cells,'cells')
6 ! maps
7 call op_decl_map(edges,nodes,2,edge ,pedge ,'pedge' )
8 call op_decl_map(edges,cells,2,ecell,pecell,'pecell')
9 ! data
10 call op_decl_dat(nodes,2,'real(8)',x,p_x,'p_x')
11 call op_decl_dat(cells,4,'real(8)',q,p_q,'p_q')
12 call op_decl_dat(cells,1,'real(8)',adt,p_adt,'p_adt')
13 call op_decl_dat(cells,4,'real(8)',res,p_res,'p_res')
14
```

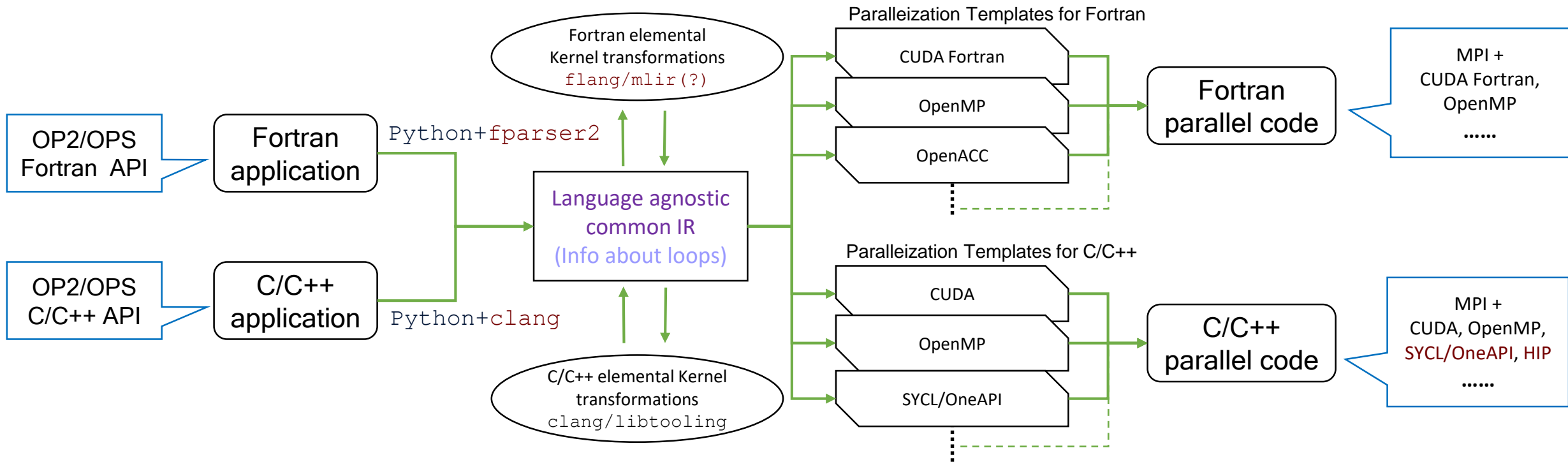


```
15 ! Elemental kernel
16 subroutine res_calc(x1,x2,q1,q2,adt1,adt2,res1,res2)
17   IMPLICIT NONE
18   REAL(kind=8), DIMENSION(2), INTENT(IN) :: x1
19   REAL(kind=8), DIMENSION(2), INTENT(IN) :: x2
20   ...
21   REAL(kind=8) :: dx,dy,mu,ri,p1,vol1,p2,vol2,f
22   dx = x1(1) - x2(1)
23   dy = x1(2) - x2(2)
24   ...
25   f = 0.5 * (vol1 * q1(1) + vol2 * q2(1)) + &
26     & mu * (q1(1) - q2(1))
27   res1(1) = res1(1) + f
28   res2(1) = res2(1) - f
29   ...
30 end subroutine
31 ! Calculate flux residual - parallel loop over edges
32 call op_par_loop_8 (res_calc, edges, &
33 & op_arg_dat(x, 1, edge, 2,"real(8)", OP_READ), &
34 & op_arg_dat(x, 2, edge, 2,"real(8)", OP_READ), &
35 & op_arg_dat(q, 1, ecell, 4,"real(8)", OP_READ), &
36 & op_arg_dat(q, 2, ecell, 4,"real(8)", OP_READ), &
37 & op_arg_dat(adt, 1, ecell, 1,"real(8)", OP_READ), &
38 & op_arg_dat(adt, 2, ecell, 1,"real(8)", OP_READ), &
39 & op_arg_dat(res, 1, ecell, 4,"real(8)", OP_INC ), &
40 & op_arg_dat(res, 2, ecell, 4,"real(8)", OP_INC ))
```

APPLICATION DEVELOPMENT



OP2/OPS CODE GENERATION



□ Simplest Code generation / translation

- Intermediate representation is simply the loop descriptions + elemental kernels
- Generated parallel code can be viewed and understood by a human !

□ Multi-layered – no opaque / black box layers

□ Built with well supported / long-term technologies - Python, Clang/libtooling, [flang?, mlir?]

EPSRC PROSPERITY PARTNERSHIP – ASIMOV [HYDRA TO OP2-HYDRA]

Virtual certification of Gas Turbine Engines

- Main consortium with partners – EPCC, Warwick, Oxford, Cambridge, Bristol and Rolls-Royce plc.

Grand Challenge 1 – Sliding Planes model of Rig250 (DLR test rig compressor)

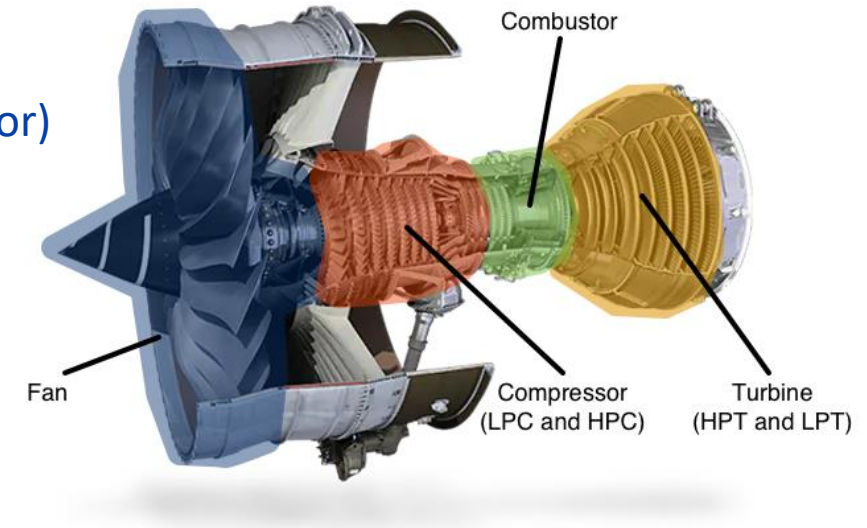
- 4.5 stage rotor-stator (10-row full annulus) | 4.58B mesh nodes.
- Need to obtain 1 revolution of compressor in less than 24 hours
- Current production estimates at 7 days

Setup

- Moving rotor-stator – sliding planes interfaces
- Rotors and Stators modelled with Hydra CFD suite – URANS (360 degree models)
- 10 rotor-stator interfaces
- Code coupling for sliding planes – move from current monolithic (Hydra only) production code to coupling

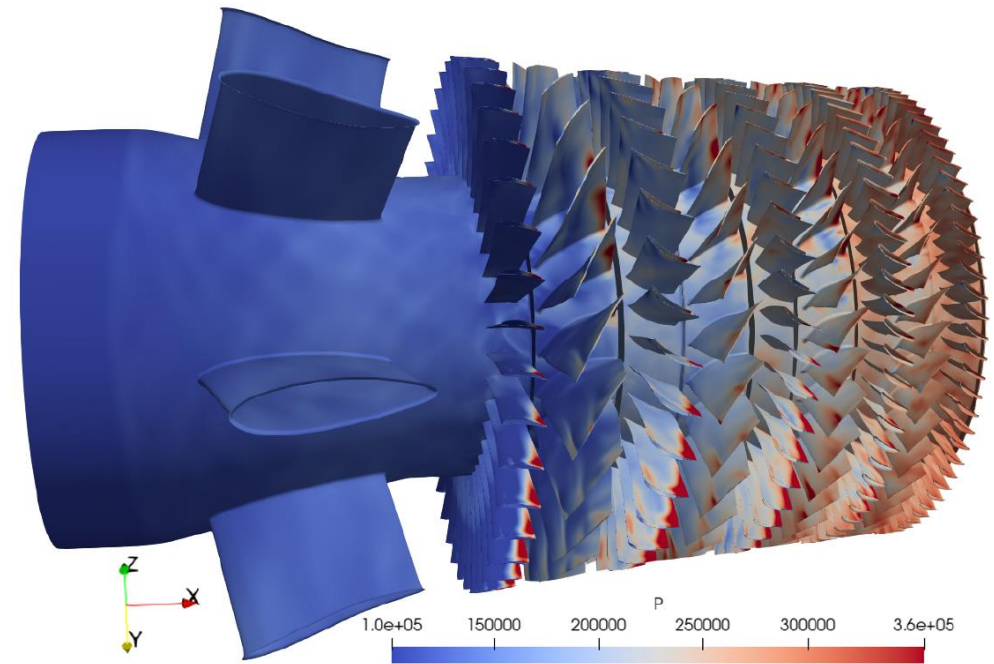
Challenges

- Performance portability – run both CPUs and GPUs by multiple vendors
- Preserve production code’s scientific code and structure – cannot re-write, MUST “evolve” not overhaul !
- Convince users to adopt ! (Ongoing for nearly 10 years now)



OP2-HYDRA PERFORMANCE

System	ARCHER2 HPE Cray EX [6]	Cirrus SGI/HPE 8600 GPU Cluster [4]
Processor	AMD EPYC 7742 @ 2.25 GHz	Intel Xeon Gold 6248 (Cascade Lake) @ 2.5 GHz + NVIDIA Tesla V100-SXM2-16GN GPU
(procs×cores) /node	2×64	2×20 + 4×GPUs
Memory/node	256 GB	384 GB + 40GB/GPU
Interconnect	HPE Cray Slingshot 2×100 Gb/s bi-directional/node	Infiniband FDR, 54.5 Gb/s
OS	HPE Cray LE (based on SLES 15)	Linux CentOS 7
Compilers	GNU 10.2.0	nvfortran (nvhpc 21.2)
Compiler Flags	-O2 -eF -fPIC	
Power/node	660W	≈ 900W



OP2-HYDRA PERFORMANCE *

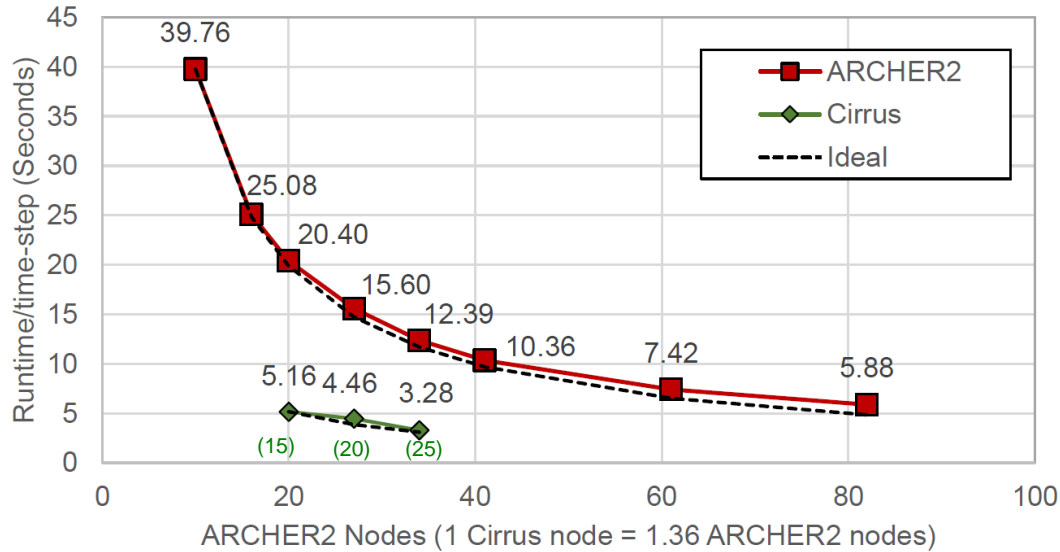
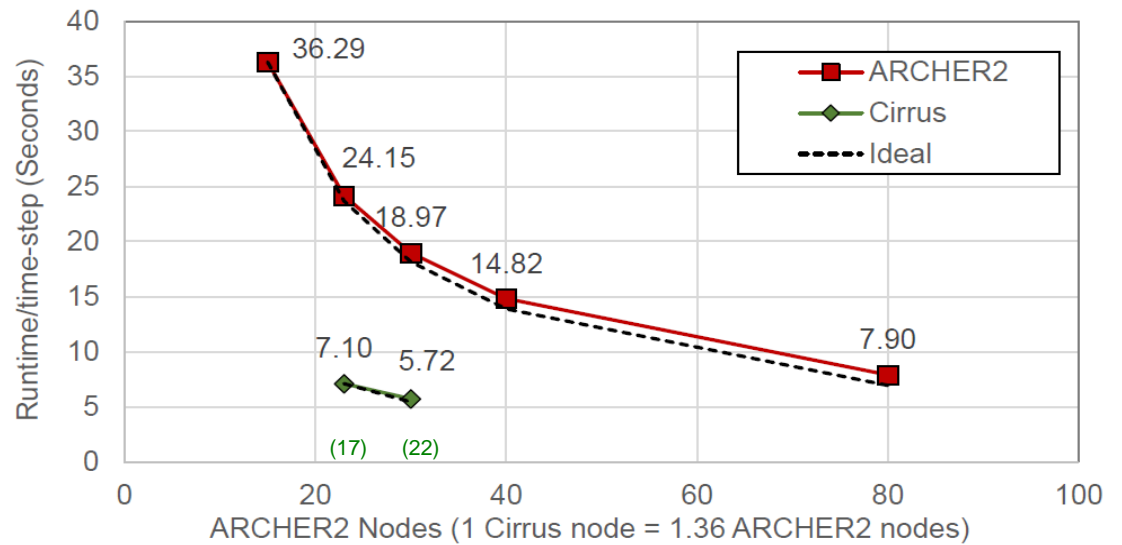


Figure 7: Rig250 1 – 10_{430M} Mesh Runtime

- **ARCHER2 @ 34 nodes**
 - 94% parallel efficiency
 - 10% coupling overhead
- **Cirrus @ 25 nodes**
 - 94% parallel efficiency
 - 20% coupling overhead
- **ARCHER2 @ 82 nodes**
 - 82% parallel efficiency
 - 20% coupling overhead

3.7- 4x speedup



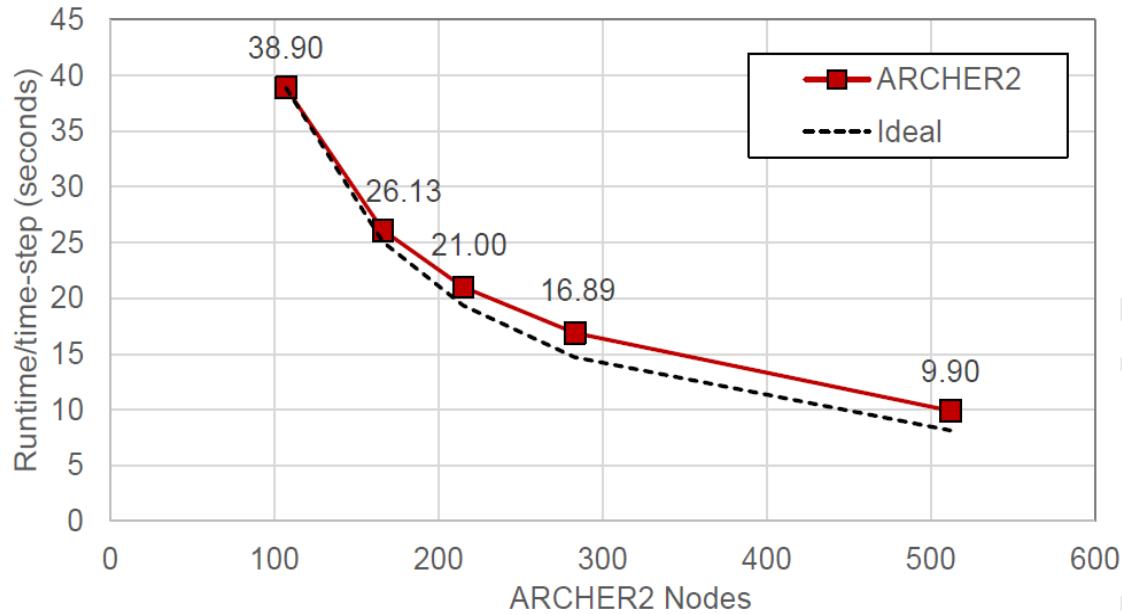
Rig250 1 – 2_{653M} Mesh Runtime

- **ARCHER2 @ 80 nodes**
 - 88% parallel efficiency
 - 8% coupling overhead
- **Cirrus @ 22 nodes**
 - 94% parallel efficiency
 - 12% coupling overhead
- **Cirrus @ 22 nodes**
 - 94% parallel efficiency
 - 12% coupling overhead

3.3 - 3.4x speedup

* Results under review

OP2-HYDRA PERFORMANCE



Rig250 1 - 10_{4.58B} Mesh Runtime

- ARCHER2 @ 512 nodes:
 - 82% parallel efficiency (vs 107 node run)
 - 15% coupling overhead

Achieved (A) and Projected (P) times to solution (hours) : Rig250, 1 revolution

Rig250 Problem	ARCHER2		Cirrus	
	Runime	#nodes	Runtime	#nodes
1 - 10 _{430M} - Monolithic	93.0 (P)	8		
1 - 10 _{430M} - Coupled	85.0 (P)	8	2.9 (P)	15
1 - 10 _{430M} - Coupled	3.3 (P)	80	1.8 (P)	25
1 - 2 _{653M} - Monolithic	110.0 (P)	8		
1 - 2 _{653M} - Coupled	40.0 (P)	8	3.9 (P)	17
1 - 2 _{653M} - Coupled	8.2 (P)	40	3.2 (P)	22
1 - 10 _{4.58B} - Coupled	14.5 (A)	166	<u>4.7 (P)</u>	122
1 - 10 _{4.58B} - Coupled	9.4 (A)	256		
1 - 10 _{4.58B} - Coupled	<u>5.5 (A)</u>	512		

- 122 Cirrus nodes is power equivalent to 166 ARCHER2 nodes
- ARCHER2 needs just over 3x more number of *power equivalent* nodes (512) to match Cirrus's runtime (4.7 hours)

* Results under review

OTHER PRODUCTION APPLICATIONS – OPENSBLI (UNI. OF SOUTHAMPTON)

Compressible Navier-Stokes solver

- With shock capturing WENO/TENO
- 4th order Finite Difference
- Single/double precision

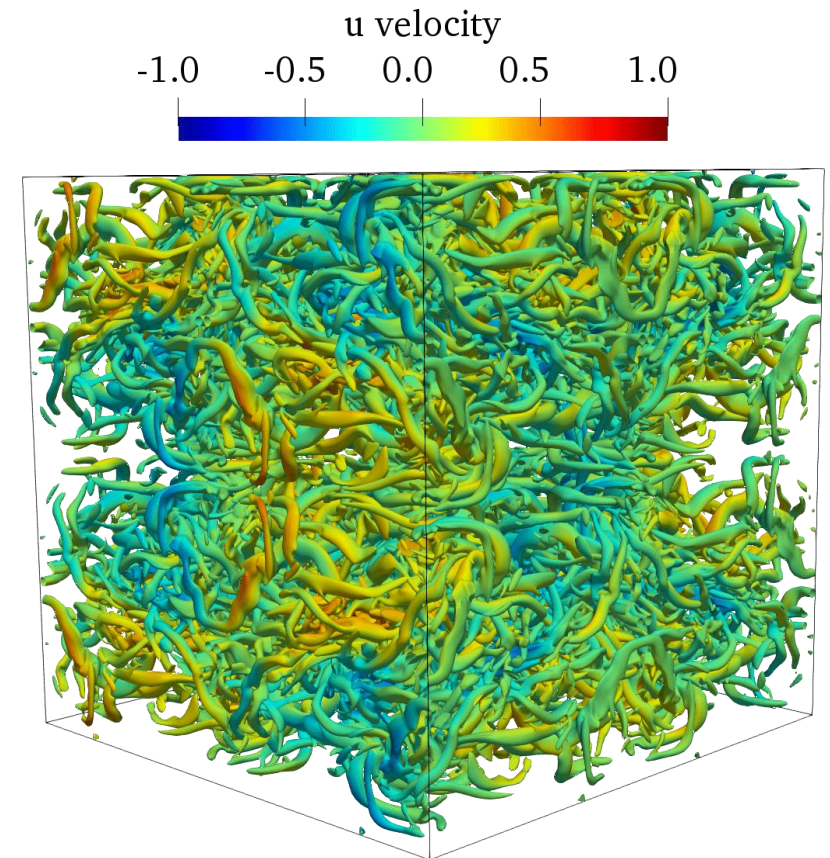
OpenSBLI is a Python framework

- Write equations in SymPy expressions
- OPS code generated

```
1 ndim = 3
2 sc1 = "**{'scheme':'Teno'}"
3 # Define the compressible Navier-Stokes equations in Einstein notation.
4 mass = "Eq(Der(rho,t), -Conservative(rhou_j,x_j,%s))" % sc1
5 momentum = "Eq(Der(rhou_i,t), -Conservative(rhou_i*u_j + KD(_i,_j)*p,x_j , %s) + Der(tau_i_j,x_j) )" % sc1
6 energy = "Eq(Der(rhoE,t), -Conservative((p+rhoE)*u_j,x_j , %s) - Der(q_j,x_j) + Der(u_i*tau_i_j ,x_j) )" % sc1
7 stress_tensor = "Eq(tau_i_j, (mu/Re)*(Der(u_i,x_j)+ Der(u_j,x_i) - (2/3)* KD(_i,_j)* Der(u_k,x_k))"
8 heat_flux = "Eq(q_j, (-mu/((gama-1)*Minf*Minf*Pr*Re))*Der(T,x_j))"
9 # Numerical scheme selection
10 Avg = RoeAverage([0, 1])
11 LLF = LLFTeno(teno_order, averaging=Avg)
12 cent = Central(4)
13 rk = RungeKuttaL5(3, formulation='SSP')
14 # Specifying boundary conditions
15 boundaries[direction][side] = IsothermalWallBC(direction, 0, wall_eqns)
16 # Generate a C code
17 alg = TraditionalAlgorithmRK(block)
18 OPSC(alg)
```

OpenSBLI

<https://opensbli.github.io/>



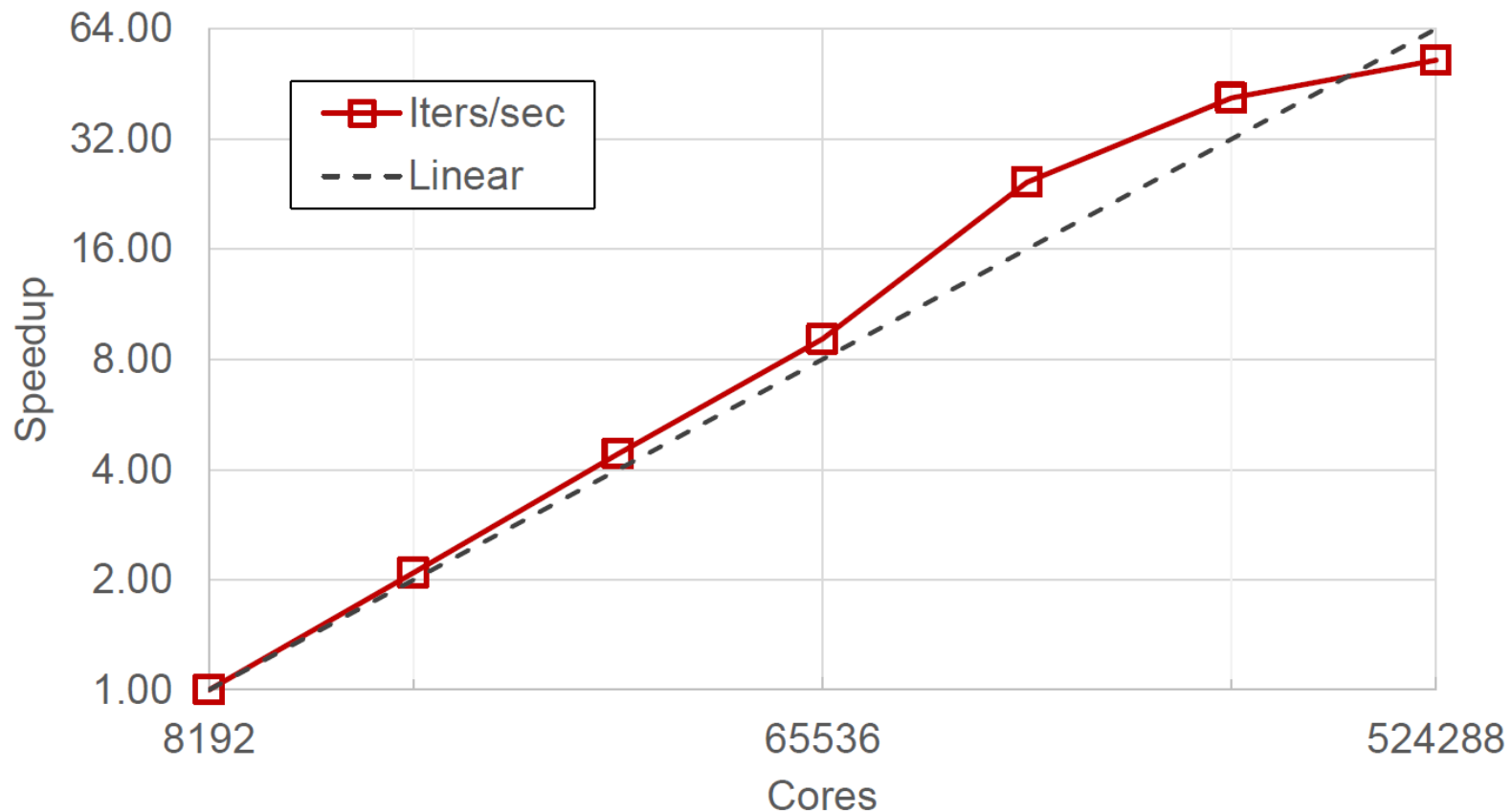
Jacobs, C. T., Jammy, S. P., Sandham N. D. (2017). OpenSBLI: A framework for the automated derivation and parallel execution of finite difference solvers on a range of computer architectures. *Journal of Computational Science*, 18:12-23, DOI: 10.1016/j.jocs.2016.11.001

OPENSBLI ON ARCHER2

□ Taylor – Green Vortex Problem – ARCHER2 benchmark

- Strong Scaling - 1024^3 Mesh
- Double precision
- Speedup calculated from 1000 iterations – includes start up time.

From recent benchmarking runs done by Andrew Turner and the ExCALIBUR Benchmarking team (Oct 2021)



OTHER PROJECTS USING OP2/OPS

- ❑ ETH Zurich – BASEMENT code (Basic Simulation Environment for Computation of Environmental Flows and Natural Hazard Simulations)
 - Flood forecast and mitigation, River morphodynamics, Design of hydraulic structures
 - Finite volume discretisation, cell centred
 - Targeting OP2 for GPU and multi-core parallelisation

- ❑ STFC – HiLeMMS project (High-Level Mesoscale Modelling System):
 - high-level abstraction layer over OPS for the solution of the Lattice Boltzmann method
 - Adaptive mesh refinement - Chombo (Lawrence Berkeley National Labs)

- ❑ University of Nottingham – CFD code development with OPS
 - Simulation of Turbomachinery flows
 - Implicit solvers using OPS's Tridiagonal Solver API

CHALLENGES – COST / EFFORT OF CONVERSION

- ❑ Converting legacy code is time consuming
 - Large code base,
 - Defunct 3rd party libs,
 - Fortran 77 or older !

- ❑ Difficult to validate code
 - New code giving the same accurate scientific output ?
 - What code should I certify ? High-level code/generated code ?
 - Difficult to convince users to use new code - fear of an opaque compiler / intermediate representation / black box !

- ❑ Incremental conversion – loop by loop
 - Simpler than CUDA, but more difficult than OpenACC/OpenMP
 - Automated conversion ?

- ❑ Changing user requirements
 - Wanting to use a DSL for doing things beyond what it was intended for !
 - Asking for “back-doors” / “escape hatches” -- leads to poor performance

CHALLENGES – CODE-GENERATION

❑ Tools not entirely mature

- Currently source-to-source with Python
- Pushing clang/LLVM source-to-source to do what we want
- What about Fortran - may be F18/Flang ?
- MLIR appearing to give some advance capabilities – see ExCALIBUR xDSL project (Tobias Grosser, Paul Kelly et al.)

❑ Code-generation for more exotic architectures – e.g. FPGAs

- Large design space
- Complex source transformations –cross loop, loop fusion and unrolling to create longer and longer pipelines !

❑ Maintainable/long term source-to-source technologies

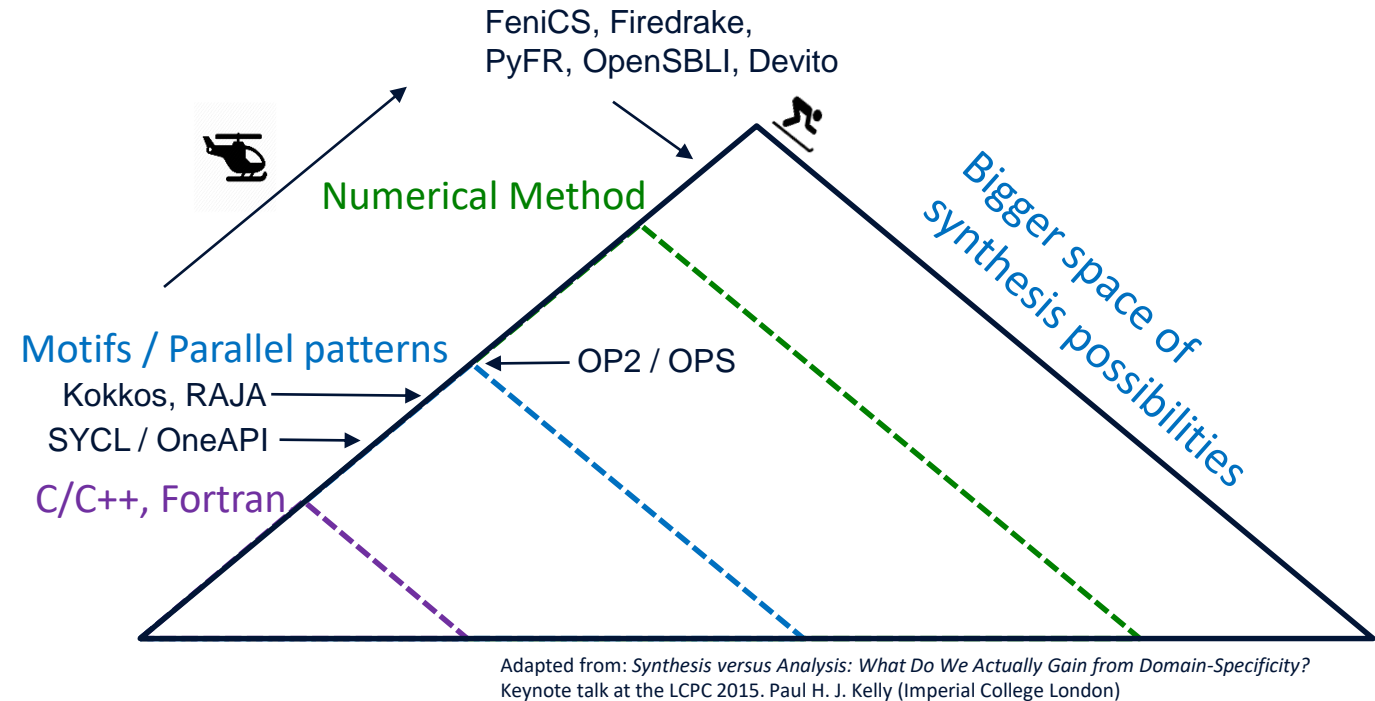
- Domain Scientists not having expertise to understand / maintain DSLs

CHALLENGES – WHO MAINTAINS THE DSL, WHAT DSL TO CHOOSE ?

- ❑ Currently purely done via academic and (small/short term) industrial funding
- ❑ Long term funding and maintenance
 - Once established probably will not be different to any other classical library
 - Will require compiler expertise to maintain code generation tools
- ❑ What DSL to choose ?
 - Re-use technologies / DSLs – especially code-gen tools (best not to reinvent !)
- ❑ Skills Gap
 - Programme in C/C++/Fortran (at a minimum)
 - Knowledge of compilers / code-generation
 - Compete for applicants – Communicate what we do better | impact of HPC / Computational Sciences
 - [In the UK] Salary ☹
 - [In the UK] Contracts ☹

DSLs / HIGH-LEVEL ABSTRACTIONS GAINING TRACTION !

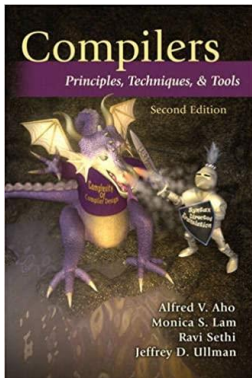
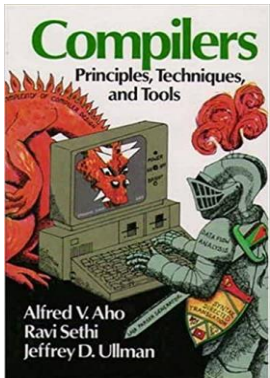
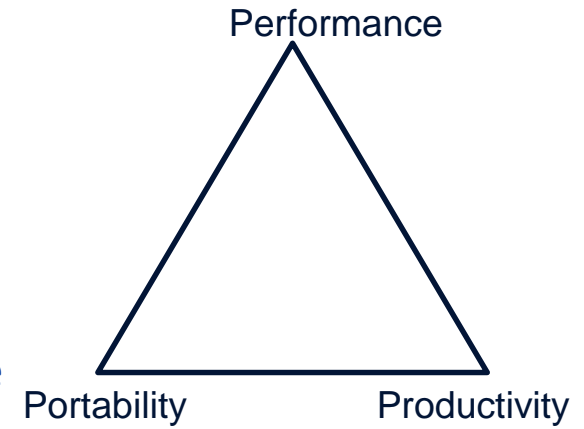
- ❑ **FEniCS** - PDE solver package - <https://fenicsproject.org/>
- ❑ **Firedrake** - automated system for the portable solution of PDEs using the finite element method
<https://www.firedrakeproject.org/>
- ❑ **PyFR** - Python based framework for solving advection-diffusion type problems on streaming architectures using the Flux Reconstruction approach - <http://www.pyfr.org/>
- ❑ **Devito** - prototype DSL and code generation framework based on SymPy for the design of highly optimised finite difference kernels for use in inversion methods -
<http://www.opesci.org/devito-public>
- ❑ **GungHO** project - Weather modelling codes (MetOffice)
- ❑ **STELLA** – DSL for stencil codes, for solving PDEs (Metro Swiss)
- ❑ **Liszt** – Stanford University : DSL for solving mesh-based PDEs -
<http://graphics.stanford.edu/hackliszt/>
- ❑ **Kokkos** – C++ template library – SNL
- ❑ **RAJA** - C++ template libraries - LLNL



Separation of Concerns – One of the four pillars of ExCALIBUR

LESSONS LEARNT AND CONCLUSIONS

- Utilizing domain knowledge will expose things that the compiler does not know
 - Iterating over the same mesh many times without change
 - Mesh is partitioned and colourable
- Compilers are conservative
 - Force it to do what you know is right for your code !
- Let go of the conventional wisdom that higher abstraction will not deliver higher performance
 - Higher abstraction leads to a bigger space of code synthesis possibilities
 - We can automatically generate significantly better code than what (most) people can (reasonably) write
 - Do not destroy performance portability by (hand-) tuning at a very low level to a specific platform



“Fundamentals and abstractions have more staying power than the technology of the moment”
Alfred Aho and Jeffrey Ullman (Turing Award Recipients 2020)

DOWNLOADS AND MORE INFORMATION

□ GitHub Repositories

- OP2 – <https://github.com/OP-DSL/OP2-Common>
- OPS – <https://github.com/OP-DSL/OPS>
- OP-DSL Webpage - <https://op-dsl.github.io/>

□ Contact

Gihan Mudalige (Warwick) - g.mudalige@warwick.ac.uk
Istvan Reguly (PPCU – Hungary) - reguly.istvan@itk.ppke.hu



ACKNOWLEDGEMENTS

- ❑ OP2 was part-funded by the UK Technology Strategy Board and Rolls-Royce plc. through the SILOET project, and the UK EPSRC projects EP/I006079/1, EP/I00677X/1 on Multi-layered Abstractions for PDEs.
- ❑ OPS was part-funded by the UK Engineering and Physical Sciences Research Council projects EP/K038494/1, EP/K038486/1, EP/K038451/1 and EP/K038567/1 on “Future-proof massively-parallel execution of multi-block applications” and EP/J010553/1 “Software for Emerging Architectures” (ASEArch) project.
- ❑ This research is supported by Rolls-Royce plc., and by the UK EPSRC (EP/S005072/1) Strategic Partnership in Computational Science for Advanced Simulation and Modelling of Engineering Systems (ASiMoV).
- ❑ Gihan Mudalige was supported by the Royal Society Industrial Fellowship Scheme (INF/R1/180012)
- ❑ Research was part-supported by the Janos Bolyai Research Scholarship of the Hungarian Academy of Sciences.
- ❑ The research has been carried out within the project Thematic Research Cooperation Establishing Innovative Informatic and Info-communication Solutions, which has been supported by the European Union and co-financed by the European Social Fund under grant number EFOP-3.6.2-16-2017-00013.
- ❑ OpenSBLI was part-funded by EPSRC grants EP/K038567/1 and EP/L000261/1, and European Commission H2020 grant 671571 “ExaFLOW: Enabling Exascale Fluid Dynamics Simulations