

Performance evaluation of explicit finite difference algorithms with varying amounts of computational and memory intensity

Satya P. Jammy^{a,*}, Christian T. Jacobs^a, Neil D. Sandham^a

^a*Aerodynamics and Flight Mechanics Group, Faculty of Engineering and the Environment, University of Southampton, University Road, Southampton, SO17 1BJ, United Kingdom*

Abstract

Future architectures designed to deliver exascale performance motivate the need for novel algorithmic changes in order to fully exploit their capabilities. In this paper, the performance of several numerical algorithms, characterised by varying degrees of memory and computational intensity, are evaluated in the context of finite difference methods for fluid dynamics problems. It is shown that, by storing some of the evaluated derivatives as single thread- or process-local variables in memory, or recomputing the derivatives on-the-fly, a speed-up of ~ 2 can be obtained compared to traditional algorithms that store all derivatives in global arrays.

Keywords: Computational Fluid Dynamics; Finite Difference Methods; Algorithms; Exascale; Parallel Computing; Performance

1. Introduction

2 Explicit finite difference methods are an important class of numerical
3 methods for the solution of partial or ordinary differential equations. For
4 example, they are used for numerically solving the governing equations in
5 computational fluid dynamics (CFD), astrophysics, seismic wave simulations,
6 financial simulations, etc.

7 In CFD they are used by many researchers for the Direct Numerical
8 Simulation (DNS) or Large Eddy Simulation (LES) of compressible flows.
9 DNS is often performed to study boundary layers, aerofoils (involving both

*Corresponding author.

E-mail address: s.p.jammy@soton.ac.uk

10 hydrodynamics and noise computations) [1], mixing analysis [2], shock-wave
11 boundary layer interactions [3] or benchmark test cases such as the Taylor-
12 Green vortex [4], decaying homogeneous isotropic turbulence, etc. Even with
13 the advances in computing hardware during the past decade, the current
14 capabilities of DNS are limited to moderate Reynolds number flows [5].

15 It is expected that computing architectures will be capable of exaFLOPs
16 (10^{18} Floating Point Operations) by 2018 and 30 exaFLOPs by 2030 [6].
17 Exascale architectures have the capability to perform DNS of the aforemen-
18 tioned examples (amongst others) at higher Reynolds numbers, or potentially
19 wall-modelled LES of the full model of an aircraft at operating Reynolds
20 numbers. However, while there is a consensus [6] that future architectures
21 would not look much like the present IBM Blue Gene, Cray, or IBM Produc-
22 tive, it is hard to predict the architectural design of such exascale systems.
23 For example, they are expected to comprise less memory per core than the
24 existing architectures. Exploiting the full potential of the exascale architec-
25 tures poses many challenges to researchers, such as the sustainability of the
26 solver's implementation with the uncertainty of architectures, the need for
27 new revolutionary algorithms/numerical methods, increasing computation
28 to communication ratio and the likelihood of I/O bottlenecks.

29 To address the problem of sustainability, taking into account the uncer-
30 tainty in future architectures, one solution adopted by the CFD community
31 involves decoupling the work of a domain scientist and a computational/-
32 computer scientist [7]. In this approach, Domain Specific Languages (DSL)
33 are developed by the computational/computer scientists, and the specifics of
34 the problem and the numerical solution method are specified in the DSL by
35 the domain scientist. Using source-to-source translation the numerical solver
36 is targetted towards different parallel hardware backends (e.g. MPI, CUDA,
37 OpenMP, OpenCL, and OpenACC) [8, 9]. This ensures that, for new archi-
38 tectures, only the backend that interfaces with the new architecture needs to
39 be written and supported by the translator. The underlying implementation
40 of the solver remains the same, thereby introducing a separation of concerns.

41 On the algorithms front, a lot of effort has gone into rewriting CFD
42 solvers to exploit the available FLOPS of existing architectures. While the
43 architectures have changed drastically in the last decade, algorithms have
44 not advanced at a similar pace [6]. Some algorithmical changes have been
45 attempted by [10, 11] to reduce the data transfer on Graphics Processing
46 Units (GPUs), but a complete and detailed study on the performance of
47 such algorithms on the existing CPU-based architectures is currently lack-
48 ing. A first step towards exascale computing would be to evaluate the per-
49 formance of algorithms characterised by varying intensities of memory usage

50 and computational cost on current CPU-based architectures for a relevant
 51 hydrodynamic test case, solved using a finite difference scheme.

52 To facilitate these investigations, the capabilities of the recently devel-
 53 oped OpenSBLI framework [12] are extended to easily generate algorithms
 54 with varying amounts of computational and memory intensity. OpenSBLI
 55 is a framework for the automated derivation and parallel execution of finite
 56 difference-based models. It is written in Python and uses SymPy to generate
 57 a symbolic representation of the governing equations and discretisation. The
 58 framework generates OPS-compliant C code that is targetted towards MPI
 59 via the OPS active library [9]. A similar approach can also be applied to any
 60 set of compute-intensive equations solved using finite difference methods.

61 The aims of this paper are to: (a) study the performance of various
 62 algorithms on current multi-core CPU-based architectures, (b) identify the
 63 best possible algorithm for the solution of explicit finite difference methods
 64 on current multi-core CPU-based architectures, and (c) demonstrate the
 65 ease at which algorithmic manipulations can be achieved with OpenSBLI
 66 framework to overcome the challenges exascale architectures can pose.

67 The rest of the paper is organised as follows. The various algorithms
 68 are described in section 2. The validation of the algorithms is presented in
 69 section 3. The performance and scaling results are presented in section 4.
 70 Some conclusions are drawn in section 5.

71 2. Algorithms

All the algorithms presented in this paper solve the three-dimensional un-
 steady compressible Navier-Stokes equations, with constant viscosity, given
 by

$$\frac{\partial \rho}{\partial t} = -\frac{\partial}{\partial x_j} [\rho u_j], \quad (1)$$

$$\frac{\partial \rho u_i}{\partial t} = -\frac{\partial}{\partial x_j} [\rho u_i u_j + p \delta_{ij} - \tau_{ij}], \quad (2)$$

and

$$\frac{\partial \rho E}{\partial t} = -\frac{\partial}{\partial x_j} [\rho E u_j + u_j p - q_j - u_i \tau_{ij}], \quad (3)$$

for the conservation of mass, momentum and energy, respectively. The quan-
 tity ρ is the fluid density, u_i is the velocity vector, p is pressure and E the
 total energy. The stress tensor τ_{ij} is defined as,

$$\tau_{ij} = \frac{1}{\text{Re}} \left(\frac{\partial u_i}{\partial x_j} + \frac{\partial u_j}{\partial x_i} - \frac{2}{3} \delta_{ij} \frac{\partial u_k}{\partial x_k} \right), \quad (4)$$

where δ_{ij} is the Kronecker Delta and Re is the Reynolds number. The heat flux term q_j is given by,

$$q_j = \frac{1}{(\gamma - 1) M^2 \text{Pr} \text{Re}} \frac{\partial T}{\partial x_j}, \quad (5)$$

where, T is temperature, M is the Mach number of the flow, Pr is Prandtl number and γ is the ratio of specific heats. The pressure and temperature are given by,

$$p = (\gamma - 1) \left(\rho E - \frac{1}{2} \rho u_j^2 \right), \quad (6)$$

and

$$T = \frac{\gamma M^2 p}{\rho}, \quad (7)$$

72 respectively. The variables that are advanced in time ($\rho, \rho u_i, \rho E$) are re-
 73 ferred to as the conservative variables, and the right-hand sides (RHS) in
 74 the mass, momentum and energy equations are referred to as the residuals
 75 of the equations.

The mass, momentum and energy equations are discretised in space using a fourth-order central finite-difference scheme and a low storage Runge-Kutta (RK) scheme with three stages of temporal discretisation. For improved stability, the convective terms in the governing equations are rewritten using the formulation of [13],

$$\frac{\partial}{\partial x_j} \rho \phi u_j = \frac{1}{2} \left(\frac{\partial}{\partial x_j} \rho \phi u_j + u_j \frac{\partial}{\partial x_j} \rho \phi + \rho \phi \frac{\partial}{\partial x_j} u_j \right), \quad (8)$$

76 where ϕ is 1, u_i or internal energy (e) for the mass, momentum and energy
 77 equations, respectively. To improve the stability of the present scheme, the
 78 viscous terms in the momentum and energy equations are expanded to second
 79 derivatives as used by [2, 10, 14].

80 A generic pseudo-code of the solution algorithm is shown in figure 2. The
 81 time loop is the most computationally expensive part of the algorithm. It
 82 consists of evaluating the primitive variables (p, u_i, T), spatial derivatives,
 83 the residual for each equation and advancing the solution in time. This
 84 is achieved by iterating over the solution points of the grid, referred to as
 85 the grid loop in the rest of the paper. Various algorithms used for the
 86 evaluation of the residual of the equations are presented herein. Starting
 87 with a memory-intensive algorithm representing a typical handwritten CFD
 88 solver, the amount of memory used and the computational intensity are
 89 varied, either by re-evaluating the derivatives on-the-fly or evaluating the

```

set-the-initial-condition
for each-iteration do
  save-state
  for each-rk-substep do
    evaluate-u_i,p,T
    evaluate-the-derivatives
    evaluate-the-residual-of-the-equations
    boundary-conditions
    advance-solution-in-time
  end for // end of rk sub loop
end for // end of iteration loop

```

Figure 1: Pseudo-code for the solution of the compressible Navier-Stokes equations.

90 derivatives using process-local variables. In all the algorithms presented, the
 91 primitive variables are evaluated and stored in memory.

92 *Baseline algorithm (BL)*. This algorithm incorporates features similar to a
 93 typical handwritten static algorithm (i.e. the derivatives in the residual of
 94 each equation are evaluated and stored in memory as arrays of grid point
 95 values; these are referred to as work arrays in the rest of the paper) on
 96 CPUs, to run as a sequential or parallel using MPI or OpenMP. For multi-
 97 threaded parallel programs, this requires the algorithm to be thread-safe in
 98 order to avoid race conditions; these occur when a variable is updated in
 99 the grid loop and the updated variable is used to update another variable
 100 in the same loop. For example, in the evaluation of the primitive variables
 101 from the conservative variables, the equation for pressure (6) is dependent
 102 on the evaluated velocity components, and the equation for temperature
 103 (7) is dependant on the evaluated pressure. When running on threaded
 104 architectures, this potentially results in race conditions. This means that
 105 temperature could be evaluated before evaluating the pressure, and pressure
 106 could be evaluated before the velocity components are evaluated. Similar
 107 candidates for race conditions exist in the update equations (which advance
 108 the conservative variables forward in time) of the RK scheme.

109 To remove the race conditions, the code is generated such that no variable
 110 is updated and used in the same loop. This is achieved by separating the
 111 evaluations into multiple loops over grid points. For example, in the evalua-
 112 tion of primitive variables, the velocity components (u_0, u_1, u_2) are grouped
 113 into a single loop as the evaluations are independent, but the pressure and
 114 temperature are evaluated in different loops.

115 When generating the code that implements the BL algorithm, the first
116 and second derivatives in the equations are evaluated and stored in work
117 arrays in order to compute the RHS residual. The evaluation of the derivative
118 of a combination of variables (e.g. $\partial(\rho u_0 u_0)/\partial x_0$) is achieved in two stages.
119 In the first stage the function $\rho u_0 u_0$ is evaluated and stored in a work array.
120 In the second stage the derivative of the work array is evaluated using the
121 central finite difference formula, and this result is stored in a new work
122 array. The work array used in the first stage is not freed in memory, but is
123 overwritten/reused when evaluating other quantities.

124 The baseline algorithm is optimised such that computationally-expensive
125 divisions are minimised. Rational numbers (e.g. finite difference stencil
126 weights) and all the negative powers of constants in the equations are evalu-
127 ated and stored at the start of the simulation. Typically, these are γ^{-1} , Pr^{-1} ,
128 Re^{-1} , and so on.

```

129 1 ndim=3 # Problem dimension
130 2 # Define the compressible Navier-Stokes in Einstein notation.
131 3 mass="Eq(Der(rho, t), - Skew(rho*u_j, x_j))"
132 4 momentum="Eq(Der(rhou_i, t), -Skew(rhou_i*u_j, x_j)-Der(p, x_i)+Der(
133     tau_i_j, x_j))"
134 5 energy="Eq(Der(rhoE, t), -Skew(rhoE*u_j, x_j)-Conservative(p*u_j,
135     x_j)+Der(q_j, x_j)+Der(u_i*tau_i_j, x_j))"
136 6 equations=[mass, momentum, energy]
137 7 # Substitutions
138 8 stress_tensor="Eq(tau_i_j, (1/Re)*(Der(u_i, x_j)+ Der(u_j, x_i)
139     -(2/3)* KD(_i, _j)* Der(u_k, x_k)))"
140 9 heat_flux="Eq(q_j, (1/((gama-1)*Minf*Minf*Pr*Re))*Der(T, x_j))"
141 0 substitutions=[stress_tensor, heat_flux]
142 1 # Define all the constants in the equations
143 2 constants=["Re", "Pr", "gama", "Minf"]
144 3 # Define coordinate direction symbol (x) this will be x_i, x_j,
145     x_k
146 4 coordinate_symbol="x"
147 5 # Formulas for the variables used in the equations
148 6 velocity="Eq(u_i, rhou_i/rho)"
149 7 pressure="Eq(p, (gama-1)*(rhoE - (1/2)*rho*(u_j*u_j)))"
150 8 temperature="Eq(T, gama*Minf*Minf/rho)"
151 9 formulas=[velocity, pressure, temperature]
152 0 # Create the problem and expand the equations.
153 1 problem=Problem(equations, substitutions, ndim, constants,
154     coordinate_symbol, metrics, formulas)
155 2 ex_eq=problem.get_expanded(problem.equations)
156 3 ex_form=problem.get_expanded(problem.formulas)
157 4 ss=Central(4) # Fourth-order central differencing
158 5 ts=RungeKutta(3) # Third-order RK scheme
159 6 # Create a numerical grid of solution points

```

```

1607 np=[64]*ndim; deltas=[2.0*pi/np[i] for i in range(len(length))]
1628 grid=Grid(ndim,{ 'delta':deltas, 'number_of_points':np})
1629 # Perform the discretisation
1630 sd=SpatialDiscretisation(ex_eq,ex_form,grid,ss)
1641 td=TemporalDiscretisation(ts,grid,constant_dt=True,sd)
1652 # Boundary conditions
1663 bc=PeriodicBoundaryCondition(grid)
1674 for dim in range(ndim):
1685     bc.apply(td.prognostic_variables,dim)
1696 # Constant initial conditions
1707 ics=["Eq(grid.work_array(rho),1.0)","Eq(grid.work_array(rhou0)
171     ,1.0)","Eq(grid.work_array(rhou1),0.0)","Eq(grid.work_array(
172     rhou2),0.0)","Eq(grid.work_array(rhoE),1.0)"]
1738 ics=GridBasedInitialisation(grid,initial_conditions)
1749 io=FileIO(td.prognostic_variables)# I/O
1750 # Simulation parameters
1761 var=['niter','Re','Pr','gama','Minf','precision','name','deltat'
177     ]
1782 values=[1000,1600,0.71,1.4,0.1,"double","test",3.385*10**-3]
1793 sp=dict(zip(vars,values))# dictionary
1804 code=OPSC(grid,sd,td,bc,ics,io,sp)#code generation

```

Listing 1: Key lines of the setup file for obtaining the BL algorithm.

181 A sample setup file used to generate an implementation of this algorithm
182 in OpenSBLI is shown in listing 1. All the algorithms presented next are
183 also optimised to reduce computationally expensive divisions. The setup file
184 for other algorithms is similar to the BL algorithm with extra attributes to
185 control the combinations of memory used and computational intensity.

Recompute All algorithm (RA). In contrast to the BL algorithm, the evaluation of pressure and temperature are first rewritten using the conservative variables,

$$p = (\gamma - 1) \left(\rho E - \frac{1}{2} \rho \left(\frac{\rho u_j}{\rho} \right)^2 \right), \quad (9)$$

and

$$T = \frac{\gamma M^2 p}{\rho} = \frac{\gamma (\gamma - 1) M^2 \left(\rho E - \frac{1}{2} \rho \left(\frac{\rho u_j}{\rho} \right)^2 \right)}{\rho}, \quad (10)$$

186 within the code to avoid race condition errors while fusing loops for the
187 evaluation of the primitive variables. Then, to evaluate the residual of the
188 equation, all the continuous spatial derivatives in the residual are replaced
189 by their respective finite difference formulas in the code generation stage.
190 This differs from the BL algorithm in that, instead of evaluating the deriva-
191 tives to work arrays and using them to compute the RHS residual, the code

```

for each-solution-point do
  double t1 = central difference formula for  $\partial u_0/\partial x_0$ 
  double t2 = central difference formula for  $\partial u_1/\partial x_1$ 
  :
  residual = t1 + t2 ...
end for

```

Figure 2: Pseudo-code for residual evaluation using SN algorithm.

192 generation process directly replaces the derivatives by their respective finite
 193 difference formulas such that they are recomputed every time.

194 This algorithm results in a code in which no work arrays are used for
 195 storing the derivatives. The memory required for this algorithm is therefore
 196 the least of all algorithms and the computational intensity is the highest of
 197 all. The control parameters to generate code for this algorithm are shown in
 198 listing 2.

```

199 1 grid = Grid(ndim, {"delta": deltas, "number_of_points": np})
200 2 grid.store_derivatives = False # Do not store derivatives

```

Listing 2: Control parameters to generate the code for RA algorithm.

201 *Store None algorithm (SN)*. This algorithm is similar to the RA algorithm.
 202 The difference is that, in the loop over grid points where the residuals are
 203 evaluated, each derivative in the RHS is evaluated to a single thread- or
 204 process-local variable. These variables are then used to update the residuals
 205 on a point-by-point basis, rather than storing all evaluations in a global-
 206 scope, grid-sized work array. To generate the code that implements this
 207 algorithm in OpenSBLI, the grid attribute `local_variables` should be set
 208 to `True` along with the control parameters given in listing 2. The pseudo-code
 209 for the residual evaluation as described here is provided in figure 2.

210 The memory footprint of this algorithm is similar to that of the RA
 211 algorithm, but is slightly less computationally-intensive. This is because, for
 212 example, if a derivative is evaluated to a process-local variable then it can
 213 be reused if that derivative appears in any of the equations more than once.

214 *Recompute Some algorithm (RS)*. In this algorithm, some of the derivatives
 215 (in this case, the first derivatives of the velocity components) are stored in
 216 work arrays and the remaining derivatives are replaced by their respective
 217 finite difference formulas in the residual. The evaluation of primitive vari-
 218 ables follows the same procedure as the RA algorithm. Listing 3 shows the

219 control parameters used to generate code for the RS algorithm. The memory
 220 usage for this algorithm is moderate, i.e. it is more than the RA algorithm
 221 but less than the BL algorithm.

```
222 1 grid = Grid(ndim, {"delta": deltas, "number_of_points": np})
223 2 grid.store_derivatives = False # Do not store derivatives
224 3 grid.derivatives_to_store = set(problem.
225   get_expanded_term_in_equations("Der(u_i,x_j)"))
```

Listing 3: Control parameters to generate the code for the RS algorithm

226 *Store Some algorithm (SS)*. This algorithm is a fusion of the RS and SN
 227 algorithms, such that the derivatives that are not stored in the RS algorithm
 228 are evaluated and stored in thread- or process-local variables as per the SN
 229 algorithm. Listing 4 shows the control parameters used to generate code
 230 for this algorithm. Compared to the SN algorithm, an additional nine work
 231 arrays would be required for the SS algorithm for the 3D cases, and an
 232 additional four work arrays for 2D cases, since the first derivatives of the
 233 velocity components are now stored.

```
234 1 grid = Grid(ndim, {"delta": deltas, "number_of_points": np})
235 2 grid.store_derivatives = False
236 3 grid.derivatives_to_store = set(problem.
237   get_expanded_term_in_equations("Der(u_i,x_j)"))
238 4 grid.local_variables = True
```

Listing 4: Control parameters in setup file to generate the code for the SS algorithm.

239 3. Validation

240 The baseline (BL) algorithm is validated for a 3D compressible Taylor-
 241 Green vortex problem, to check the correctness of the solver. The initial
 242 conditions and the post-processing procedure are described in [4]. The simu-
 243 lations are performed in a cube of non-dimensional length 2π , with periodic
 244 boundary conditions in all three directions for grids containing 64^3 , 128^3 , 256^3
 245 and 512^3 solution points. The Mach number, Prandtl number and Reynolds
 246 number of the flow are taken as 0.1, 0.71 and 1600, respectively. The non-
 247 dimensional time-step for the 64^3 grid size was set to 3.385×10^{-3} , and was
 248 halved for each increase in the grid size by a factor of 2^3 . Double-precision
 249 is used throughout all simulations presented in this paper.

250 Figure 3 shows the evolution of kinetic energy and enstrophy compared
 251 with the reference data [15]. The results from the BL algorithm agree very
 252 well with the reference data for the 512^3 case. For computational expedi-
 253 ence, the other algorithms are validated on the 128^3 grid. For each one, the

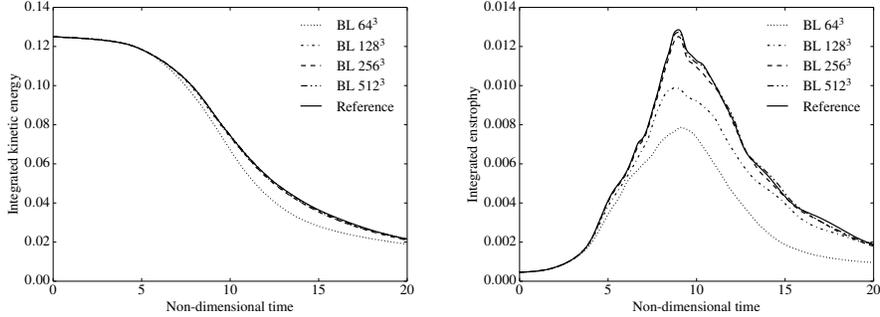


Figure 3: Left: Evolution of the integral of kinetic energy. Right: Evolution of the integral of enstrophy.

254 results relative to the BL algorithm are found to be the same up to machine
 255 precision.

256 4. Performance evaluation

257 After checking that the results from the various algorithms match, the
 258 performance of the different algorithms were evaluated using the same Taylor-
 259 Green vortex test case described in section 3. All simulations are performed
 260 on ARCHER (the UK National Supercomputing Service) and the code that
 261 implements the various algorithms is compiled using the Cray C compiler
 262 (version 2.4.2) with the -O3 optimisation flag. Each ARCHER node com-
 263 prises 24 cores, with each MPI process being mapped to its own individual
 264 core. All simulations for performance evaluation purposes are run in par-
 265 allel using 24 MPI processes/cores (one ARCHER node). The run-time of
 266 the time iteration loop was recorded for 500 iterations and is summarised in
 267 table 1 for a range of grid sizes.

N_x	N_y	N_z	BL	RA	RS	SN	SS
64	64	64	16.21	9.29	10.76	8.44	9.78
128	128	128	182.55	98.18	97.36	90.72	88.95
256	256	256	1561.52	765.42	802.76	693.66	685.25

Table 1: Total run-time in seconds for different grid sizes for all algorithms on ARCHER using 24 MPI processes.

268 The data in table 1 is plotted in figure 4; from this figure it can be inferred
 269 that when the amount of memory access is reduced, the current CPU-based

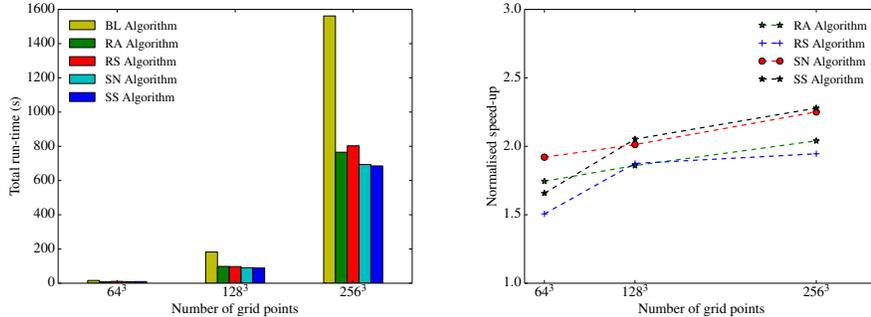


Figure 4: Left: Figure showing data in table 1. Right: speed-up of algorithms normalised with the BL algorithm.

270 architectures perform better, even though the computational intensity of
 271 such algorithms is higher. The baseline algorithm is a factor of ~ 2 slower
 272 than all the other algorithms presented in this paper. For larger grid sizes
 273 the benefit of the SS algorithm becomes more pronounced.

274 4.1. Scaling

275 Strong scaling tests were performed for the best performing algorithm
 276 (i.e. the SS algorithm) on ARCHER for the test problem with a total of
 277 1.07×10^9 grid points and the runtime of the time iteration loop was recorded
 278 for 10 iterations. Figure 5 shows the strong scaling results on ARCHER up
 279 to 73,728 MPI processes/cores (i.e. 3,072 ARCHER nodes). The minimum
 280 number of processes required for running the problem is 120. The algorithm
 281 shows a near-linear scaling (speed-up of 2) until 36,864 MPI processes (i.e.
 282 1,536 ARCHER nodes) and thereafter the speed-up is ~ 1.5 as the process
 283 count doubles.

284 Weak scaling tests were also performed for the SS algorithm. Here, the
 285 number of MPI processes was varied from 192 to 65,856 (i.e. 8 to 2,744
 286 ARCHER nodes), while the number of grid points per MPI process was kept
 287 fixed at 64^3 and the runtime of the time iteration loop was recorded for 10
 288 iterations. The largest grid size considered comprises ~ 17 billion solution
 289 points. Figure 6 demonstrates that the normalised run-time is near-ideal.

290 5. Conclusion

291 In this paper the automated code generation capabilities of the OpenS-
 292 BLI framework have been extended to easily modify the memory usage and

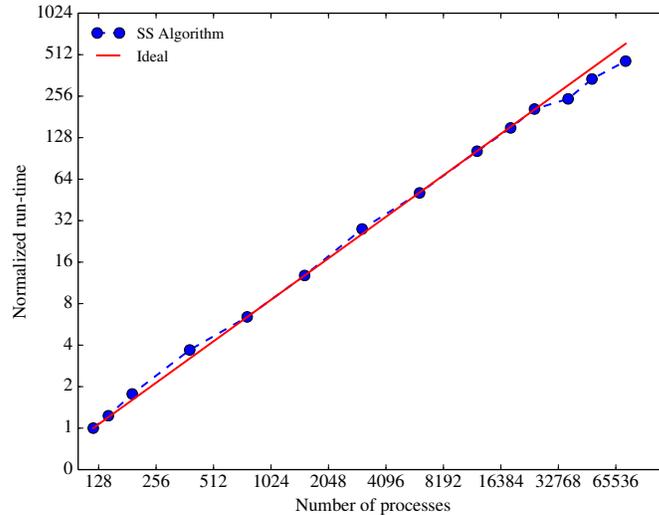


Figure 5: Strong scaling of the SS algorithm on ARCHER up to 73,728 cores using 1.07×10^9 grid points. The run-time has been normalised by that of the 120-process case.

293 computational intensity of the solution algorithm. It was found that the
 294 baseline (BL) algorithm featured in traditional CFD codes, in which all
 295 derivatives are evaluated and stored in work arrays, is not the best algorithm
 296 in terms of performance on current multi-core CPU-based architectures. Re-
 297 computing all or some of the derivatives performs better than the baseline
 298 algorithm. The best algorithm found here for the solution of the compress-
 299 ible Navier-Stokes equations is to store only the first derivatives of velocity
 300 components in work arrays, and compute the remaining spatial derivatives
 301 and store them in thread- or process-local variables. The run-time of such
 302 an algorithm has been shown to be ~ 2 times smaller than the BL algorithm.
 303 Through the use of modern code generation techniques in the OpenSBLI
 304 framework, it has been demonstrated that by changing just a few attributes
 305 (three in this case) in the problem setup file, different algorithms with vary-
 306 ing degrees of memory and computational intensity can be readily generated
 307 automatically. The methodology presented in this paper can also be used to
 308 find the best possible algorithm for other existing architectures such as GPUs
 309 or Intel Xeon Phi coprocessors. Moreover, existing numerical models that use
 310 finite difference methods for the solution of any governing equations can be
 311 optimised for the current CPU-based architectures. When exascale systems

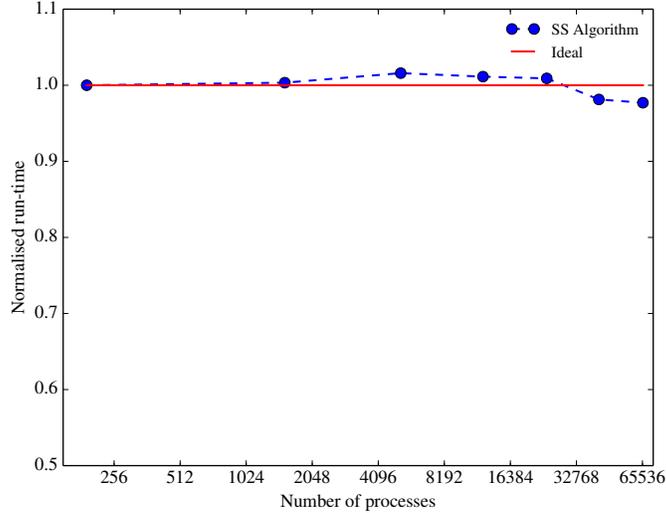


Figure 6: Weak scaling of the SS algorithm on ARCHER with 64^3 grid points per MPI process up to 65,856. The results have been normalised by the run-time from the 192-process case.

312 become available, depending on their architecture and amount of available
 313 memory, users can readily tune the memory and computational intensity in
 314 the OpenSBLI framework to determine the best performing algorithm on
 315 such systems.

316 6. Acknowledgements

317 SPJ was supported by an EPSRC grant entitled “Future-proof massively-
 318 parallel execution of multi-block applications” (EP/K038567/1). CTJ was
 319 supported by a European Union Horizon 2020 project grant entitled “Ex-
 320 aFLOW: Enabling Exascale Fluid Dynamics Simulations” (grant reference
 321 671571). The data behind the results presented in this paper will be avail-
 322 able from the University of Southampton’s institutional repository. The
 323 authors acknowledge the use of the UK National Supercomputing Service
 324 (ARCHER), with computing time provided by the UK Turbulence Consor-
 325 tium (EPSRC grant EP/L000261/1).

326 **References**

- 327 [1] L. E. Jones, R. D. Sandberg, Acoustic and hydrodynamic analysis of the
328 flow around an aerofoil with trailing-edge serrations, *Journal of Fluid*
329 *Mechanics* 706 (2012) 295–322. doi:10.1017/jfm.2012.254.
- 330 [2] S. Pirozzoli, M. Bernardini, S. Marié, F. Grasso, Early evolution of the
331 compressible mixing layer issued from two turbulent streams, *Journal of*
332 *Fluid Mechanics* 777 (2015) 196–218. doi:10.1017/jfm.2015.363.
- 333 [3] B. Wang, N. D. Sandham, Z. Hu, W. Liu, Numerical study of
334 oblique shock-wave/boundary-layer interaction considering sidewall ef-
335 fects, *Journal of Fluid Mechanics* 767 (2015) 526–561. doi:10.1017/
336 jfm.2015.58.
- 337 [4] J. DeBonis, Solutions of the Taylor-Green Vortex Problem Using High-
338 Resolution Explicit Finite Difference Methods, 51st AIAA Aerospace
339 Sciences Meeting including the New Horizons Forum and Aerospace
340 Exposition (February) (2013) 1–9. doi:10.2514/6.2013-382.
- 341 [5] W. C. Reynolds, Whither Turbulence? Turbulence at the Crossroads:
342 Proceedings of a Workshop Held at Cornell University, Ithaca, NY,
343 March 22–24, 1989, Springer Berlin Heidelberg, Berlin, Heidelberg,
344 1990, Ch. The potential and limitations of direct and large eddy simu-
345 lations, pp. 313–343. doi:10.1007/3-540-52535-1_52.
- 346 [6] J. Slotnick, A. Khodadoust, J. Alonso, D. Darmofal, W. Gropp,
347 E. Lurie, D. Mavriplis, CFD Vision 2030 Study: A Path to Revolution-
348 ary Computational Aerosciences, Nasa Cr-2014-21878 (March) (2014)
349 1–73.
- 350 [7] C. T. Jacobs, M. D. Piggott, Firedrake-Fluids v0.1: numerical mod-
351 elling of shallow water flows using an automated solution framework,
352 *Geoscientific Model Development* 8 (3) (2015) 533–547. doi:10.5194/
353 gmd-8-533-2015.
- 354 [8] F. Rathgeber, G. R. Markall, L. Mitchell, N. Lorient, D. A. Ham,
355 C. Bertolli, P. H. Kelly, PyOP2: A High-Level Framework for
356 Performance-Portable Simulations on Unstructured Meshes, in: *High*
357 *Performance Computing, Networking Storage and Analysis, SC Com-*
358 *panion*, IEEE Computer Society, 2012, pp. 1116–1123.

- 359 [9] I. Z. Reguly, G. R. Mudalige, M. B. Giles, D. Curran, S. McIntosh-
360 Smith, The OPS Domain Specific Abstraction for Multi-Block Struc-
361 tured Grid Computations, in: Proceedings of the 2014 Fourth Interna-
362 tional Workshop on Domain-Specific Languages and High-Level Frame-
363 works for High Performance Computing, IEEE Computer Society, 2014,
364 pp. 58–67. doi:10.1109/WOLFHPC.2014.7.
- 365 [10] F. Salvatore, M. Bernardini, M. Botti, GPU accelerated flow solver
366 for direct numerical simulation of turbulent flows, Journal of Computa-
367 tional Physics 235 (2013) 129–142. doi:10.1016/j.jcp.2012.10.012.
- 368 [11] J. C. Thibault, Implementation of a Cartesian Grid Incompressible
369 Navier-Stokes Solver on Multi-Gpu Desktop Platforms Using Cuda, New
370 Horizons (May) (2009) 1–15.
- 371 [12] C. T. Jacobs, S. P. Jammy, N. D. Sandham, OpenSBLI: A framework
372 for the automated derivation and parallel execution of finite difference
373 solvers on a range of computer architectures, Journal of Computational
374 Science, SubmittedarXiv:1609.01277.
- 375 [13] G. A. Blaisdell, E. T. Spyropoulos, J. H. Qin, The effect of the formu-
376 lation of nonlinear terms on aliasing errors in spectral methods, Ap-
377 plied Numerical Mathematics 21 (3) (1996) 207–219. doi:10.1016/
378 0168-9274(96)00005-0.
- 379 [14] N. Sandham, Q. Li, H. Yee, Entropy Splitting for High-Order Numer-
380 ical Simulation of Compressible Turbulence, Journal of Computational
381 Physics 178 (2) (2002) 307–322. doi:10.1006/jcph.2002.7022.
- 382 [15] Z. Wang, K. Fidkowski, R. Abgrall, F. Bassi, D. Caraeni, A. Cary,
383 H. Deconinck, R. Hartmann, K. Hillewaert, H. Huynh, N. Kroll, G. May,
384 P.-O. Persson, B. van Leer, M. Visbal, High-order cfd methods: current
385 status and perspective, International Journal for Numerical Methods in
386 Fluids 72 (8) (2013) 811–845. doi:10.1002/flid.3767.