



Citation for published version:

Mueller, E, Scheichl, R & Vainikko, E 2015, 'Petascale solvers for anisotropic PDEs in atmospheric modelling on GPU clusters', *Parallel Computing*, vol. 50, pp. 53-69. <https://doi.org/10.1016/j.parco.2015.10.007>

DOI:

[10.1016/j.parco.2015.10.007](https://doi.org/10.1016/j.parco.2015.10.007)

Publication date:

2015

Document Version

Peer reviewed version

[Link to publication](#)

Publisher Rights

CC BY-NC-ND

University of Bath

Alternative formats

If you require this document in an alternative format, please contact:
openaccess@bath.ac.uk

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

Petascale solvers for anisotropic PDEs in atmospheric modelling on GPU clusters

Eike Hermann Müller^{a,*}, Robert Scheichl^a, Eero Vainikko^b

^a*Department of Mathematical Sciences, University of Bath, BA2 7AY, Bath, United Kingdom*

^b*Institute of Computer Science, University of Tartu, Liivi 2, Tartu 50409, Estonia*

Abstract

Memory bound applications such as solvers for large sparse systems of equations remain a challenge for GPUs. Fast solvers should be based on numerically efficient algorithms and implemented such that global memory access is minimised. To solve systems with trillions ($\mathcal{O}(10^{12})$) unknowns the code has to make efficient use of several million individual processor cores on large GPU clusters.

We describe the multi-GPU implementation of two algorithmically optimal iterative solvers for anisotropic PDEs which are encountered in (semi-) implicit time stepping procedures in atmospheric modelling. In this application the condition number is large but independent of the grid resolution and both methods are asymptotically optimal, albeit with different absolute performance. In particular, an important constant in the discretisation is the CFL number; only the multigrid solver is robust to changes in this constant. We parallelise the solvers and adapt them to the specific features of GPU architectures, paying particular attention to efficient global memory access. We achieve a performance of up to 0.78 PFLOPs when solving an equation with $0.55 \cdot 10^{12}$ unknowns on 16384 GPUs; this corresponds to about 3% of the theoretical peak performance of the machine and we use more than 40% of the peak memory bandwidth with a Conjugate Gradient (CG) solver. Although the other solver, a geometric multigrid algorithm, has a slightly worse performance in terms of FLOPs per second, overall it is faster as it needs less iterations to converge; the multigrid algorithm can solve a linear PDE with half a trillion unknowns in about one second.

Keywords: iterative solver, multigrid, Graphics Processing Unit, massively parallel, atmospheric modelling

1. Introduction

Many problems in geophysical modelling require the fast solution of anisotropic partial differential equations (PDEs) in “flat” domains. For example, a global PDE for the pressure correction has to be solved in every time step of many numerical weather- and climate prediction models if an implicit method is used to advance the atmospheric fields forward in time. In this work we consider a model problem for this PDE which can be written schematically as

$$-\omega(h)^2 \left(\nabla_S^2 u + \lambda^2 \frac{1}{r^2} \frac{\partial}{\partial r} \left(r^2 \frac{\partial u}{\partial r} \right) \right) + u = f; \quad (1)$$

where $\omega(h) = \mathcal{O}(h^2)$, r is the radial coordinate and ∇_S the component of the derivative which is tangential to the surface of the earth. The exact form of the parameters ω and λ is described below. As the height of the atmosphere is much smaller than the horizontal extent of the domain in global models, after discretisation this equation has a very strong vertical anisotropy. The discretised PDE can be written as a sparse system of equations

$$\underline{A} \mathbf{u} = \mathbf{f}, \quad (2)$$

where \underline{A} is a sparse $n \times n$ matrix and the vector \mathbf{u} has n entries and represents the global (Exner-) pressure correction field in the whole domain. In modern applications the number of degrees of freedom n can be very large, for a global grid with a horizontal resolution of 1km or less and $\mathcal{O}(100 - 200)$ vertical levels a system with $n \gtrsim 10^{11}$ unknowns needs to be solved. Resolutions of this order are expected to be achieved by state-of-the-art global forecast models within the next decade. As the PDE is solved in every model time step and accounts for a significant amount of the total model runtime, it is crucial to solve it as fast as possible to deliver forecasts on operational timescales. This can only be achieved by using algorithmically optimal methods and implementing them on the fastest available hardware. The main challenge faced by the solver is the vertical anisotropy which prevents a standard approach such as a geometric multigrid algorithm with point

*Corresponding author

Email address: e.mueller@bath.ac.uk (Eike Hermann Müller)

smoother or an iterative solver with Jacobi preconditioner. Recently we have shown that Krylov subspace solvers and multigrid methods tailored to the structure of the problem are highly efficient and scale up to tens of thousands of CPU cores [1]. In particular we found that for our application the tensor-product geometric multigrid solver, suggested and first analysed in [2], is superior to a preconditioned Conjugate Gradient method and to state of the art parallel algebraic multigrid (AMG) implementations from the DUNE [3, 4] and Hypre [5] libraries.

As explained in detail below, the PDE for the atmospheric pressure correction has the structure of a shifted Laplace equation, which is usually known as the (sign-positive) Helmholtz equation in the meteorological literature. In contrast to the conventional Helmholtz equation encountered for example in wave scattering problems, the elliptic operator we consider is positive definite. While implicit time marching schemes permit larger model time steps, advective time scales and constraints on the accuracy of the solution limit the permitted time step size. This implies that the CFL number, which is proportional to the ratio of the time step size and horizontal grid spacing, is typically in the range 2 – 10. After preconditioning, the condition number of the elliptic operator is $\mathcal{O}(100 - 1000)$, independent of the horizontal grid resolution. Because of this, Krylov subspace methods and multigrid algorithms with a fixed number of levels are asymptotically optimal and algorithmically scalable as the problem size increases. However, the performance of Krylov methods deteriorates with growing CFL number.

Graphics Processing Units (GPUs) have been used very successfully in many areas of Scientific Computing and can be superior to more traditional CPU architectures both in terms of speed and power efficiency. A particular challenge for solvers of sparse systems of linear equations such as (2) is that their performance is typically limited by the speed with which data can be read from (and written to) global GPU memory. While the number of floating point operations for the iterative solvers we consider is typically two to five times larger than the number of memory operations, on modern GPUs, such as the Kepler GK110 on the K20X cards on the Titan supercomputer [6] the cost for one (double precision) memory access is more than $40\times$ larger than the cost of one floating point operation. This factor is given by the ratio of the peak floating point performance and the peak global memory bandwidth, namely $1.31 \text{ TFLOPs} / (250 \text{ GByte/s}) \times 8 \text{ Byte} \approx 42$ for double precision arithmetic on the K20X card [7]. It should be compared to the corresponding number for CPU architectures where only around 3 floating point operations can be carried out per variable loaded from memory. Furthermore, due to limited memory, only problems with up to a few million degrees of freedom can be solved on a single GPU. To solve larger systems a distributed-memory multi-GPU implementation has to be used. For problems with up to a trillion (10^{12}) unknowns, several million processor cores are necessary.

To implement the fastest possible massively parallel GPU solver we followed three design principles:

1. **Algorithmically optimal solver.** To minimise the overall solution time, the biggest gains can be achieved by using an iterative solver method which is tailored to the problem to be solved and converges in the smallest possible number of iterations. Krylov subspace methods are very popular in meteorological applications because of their simplicity (see e.g. [8, 9, 10, 11] and the detailed review in [1]). For anisotropic PDEs it is particularly important to exploit the strong coupling in the vertical direction by using a suitable preconditioner. Since the elliptic system considered in this work is symmetric and positive definite, the most suitable Krylov subspace method is a Conjugate Gradient solver preconditioned with vertical line relaxation. The preconditioner requires the frequent solution of a tridiagonal system in each vertical column; this can be achieved with the Thomas algorithm (see e.g. [12]). However, we already found in [1] that the geometric tensor-product multigrid solver proposed and analysed in [2], which uses vertical line relaxation as the smoother, converges significantly faster than the preconditioned CG iteration. As the numerical experiments in this article confirm, the CG solver requires at the order of 60 iterations to reduce the residual by five orders of magnitude, whereas the multigrid method converges in less than 10 iterations in atmospheric simulations. The situation turns even more in favour of the multigrid method if the CFL number is not restricted by accuracy constraints or the explicit treatment of other processes, such as advection.
2. **Memory optimised CUDA-C implementation.** We optimised the single GPU implementation by minimising the number of memory references per iteration. As already shown in [13], the biggest gains can be achieved by using a “matrix-free” approach and recomputing the (sparse) matrix \underline{A} instead of storing it explicitly. For example, carrying out a sparse matrix-vector product (SpMV) $\mathbf{y} \leftarrow \underline{A}\mathbf{x}$ requires 1 to 7 global reads (depending on the caching of the vector \mathbf{x}) and 1 global write at each grid cell. This should be compared to a matrix-explicit implementation which requires 7 additional reads for a finite volume stencil and can hence be more than twice as expensive. Not storing the matrix explicitly also reduces the memory requirements of the solver significantly. This allows the solution of larger problems and better utilisation of the GPU resources. For optimal global memory throughput on the GPU it is crucial to adapt the data layout to achieve optimally coalesced access for all threads in a warp. This requires a horizontally contiguous ordering of the degrees of freedom, which differs from the vertically contiguous ordering which allows optimal cache reuse on CPUs. In addition we reduced the number of memory references by fusing several GPU kernels.

We find that on a single GPU our CG implementation achieves 36% – 56% of the peak global memory bandwidth

depending on the problem size. For the multigrid solver the rate is slightly lower with 15% – 36% of the peak global memory bandwidth, but this is more than compensated by the faster convergence rate.

- 3. Massively parallel multi-GPU code.** We extended our implementation to clusters of GPUs by using a horizontal decomposition of the computational domain, which is common for applications in atmospheric modelling. For this we used the Generic Communication Library [14] which allows the straightforward implementation of halo exchanges on structured two- and three-dimensional grids and supports GPUDirect data transfer between different GPUs. For the largest problem we studied, the additional overhead from the MPI communications is about 10% for the CG solver and about 40% for the multigrid solver, both solvers show very good weak scaling to up to 16384 GPUs.

Main achievements. In this paper we describe this approach in detail for the solution of a model equation which captures the main features of the PDE for the pressure correction in global weather- and forecast models; further details on the model equation and relevant meteorological literature can be found in previous publications [1, 15]. In [13] we described the single-GPU implementation of a matrix-free CG solver, here we extend this approach to the geometric-multigrid solver analysed in [2] and extend both solvers to run on clusters of GPUs. We tested the performance of our solvers and ran them on up to 16384 GPUs of the Titan Cray XK7 cluster (OLCF, Oak Ridge National Lab), which contains 18,688 nVidia K20X cards with GK110 Kepler GPUs and is currently ranked as the second fastest computer in the world (top500.org, June 2014 [6]).

We are able to solve a problem with half a trillion ($0.55 \cdot 10^{12}$) degrees of freedom in about one second with the multigrid solver. The GPU implementation is about a factor four faster on one K20X GPU card on Titan than our optimised Fortran 90 CPU code running on one 16 core AMD Opteron processor of HECToR, the UK’s national supercomputing resource.

On Titan we achieve a performance of 0.78 PFLOPs for the CG solver (and 0.65 PFLOPs for the multigrid algorithm). It should be stressed that this time includes all components of the solver, such as host-device data transfer and transposition of the fields for horizontally contiguous ordering on the GPU. As the code is bandwidth limited, the absolute performance should be quantified in fractions of the peak global GPU bandwidth. For the CG solver we can achieve a percentage of 32%-42% of the peak global memory bandwidth when running on 16384 GPUs, for the multigrid solver this fraction is 15%-25%. All source code is freely available for download under the LGPL 3 license.

Previous work. Early work on the GPU parallelisation of Conjugate Gradient methods and multigrid solvers for sparse linear systems is discussed in [16, 17]. The authors solve the two dimensional shifted Laplace equation $-\Delta u + \sigma u = RHS$ [16] and the Poisson equation $-\Delta u = 0$ [17] arising in implicit time stepping methods for the solution of the Navier Stokes equations. In both cases the code is implemented by using the low level graphics API on GPU hardware which is quite dated now. More recent GPU implementations of (preconditioned) Conjugate Gradient- [18, 19, 20, 21, 22, 23, 24] and geometric multigrid solvers [25, 26, 27] are also reported in the literature. To solve problems which arise for example from finite element discretisations on unstructured grids, typically the system matrix is stored explicitly in formats such as compressed sparse row storage (CSR) or in the ELLPACK format (see [28] for a detailed discussion of sparse matrix storage formats on GPUs) and the authors concentrate on optimising the sparse matrix-vector multiplication. The advantage of this approach is that the solvers can be applied to very general and grid-independent problems such as power grid simulations in [25] or arbitrary sparse matrices from the University of Florida sparse matrix collection [29] as described in [19]. However, for very general problems, the construction of a suitable preconditioner is very difficult and convergence is slow. Similarly, while AMG implementations such as those reported in [30, 31] can be used to solve a very large class of elliptic problems, the requirement of explicit matrix storage makes them more expensive than geometric multigrid for the structured PDEs which we discuss in this article. On the other hand, the only matrix-free implementations we are aware of are [17, 18, 20], and in all cases the authors focus on solving the homogeneous and isotropic Poisson equation in a regular two- or three- dimensional domain. Both extremes should be compared to our approach: by exploiting the structure of the problem to construct a suitable preconditioner our solvers can deal with three dimensional anisotropic equations on curved domains but avoid explicit storage of the matrix which has a negative impact on performance for bandwidth limited applications.

More recently iterative solvers have also been parallelised across multiple GPUs. For example (unpreconditioned) Conjugate Gradient solvers are tested for a range of sparse matrices in [19, 21, 22]. Of particular interest for this work are the multigrid solvers discussed in [24, 23, 27] for the 3D Poisson equation which arises in implicit time stepping methods for the solution of the Navier Stokes equations. The equation we consider in the following can be derived in a similar fashion for the compressible Euler equations. A significant but important difference is that compressibility gives rise to an additional zero order term which introduces an intrinsic length scale beyond which interactions between gridpoints are exponentially suppressed. This has an important impact on the parallelisation of the multigrid solver: it is sufficient to use a relatively small number of multigrid levels and one or two smoother iterations are sufficient to solve the well conditioned coarse grid problem, thus avoiding an exact global coarse grid solve across all processors. This should be compared to

the more complicated approaches such as parallel coarse grid aggregation discussed in [27] for the homogeneous Poisson equation.

An interesting approach combining the computational power of both the CPU and the GPU on a node is described in [32] where the authors describe a BiCG solver for the Poisson problem on an unstructured grid and parallelise the solver on a cluster with up to 17 nodes. A multigrid V-cycle is used for preconditioning and the smoother, which is a separate multigrid iteration on structured subgrids, is offloaded to the GPU. As CPU-GPU data transfer is expensive and there are now efficient ways for exchanging data directly between GPUs on different nodes, we implemented our solver such that all calculations are carried out on the device only.

Only recently clusters with several thousands of GPUs have become available and as far as we are aware to date there are no multi-GPU implementations which have been shown to scale to up to more than around 100 GPUs. Parallel scaling on up to 128 GPU for a CG solver with Jacobi preconditioner for the Poisson equation is described in [24], and scaling for the multigrid solver of the same problem is reported on up to 64 nodes in [27]. Elliptic problems solved so far typically have less than 1 billion unknowns and the results in this work represent a significant contribution to extending the scalability of iterative solvers to several millions of processor cores on tens of thousands of GPUs and for solving very large systems with up to half a trillion unknowns. In this context we mention the work reported in [33] where the authors achieved a GPU performance of 1.9 PFLOPs for an explicit time stepping solver of the shallow water equations on a cubed sphere grid (however, in contrast to implicit timestepping methods, this does not require the solution of an elliptic PDE for the pressure correction). That solver is run on 3750 nodes with 1 GPU each to solve problems with up to 4 billion unknowns per atmospheric variable.

Although we believe that here we describe the first massively parallel GPU implementation of solvers for sparse systems with more than half a trillion ($0.5 \cdot 10^{12}$) unknowns, problems of this size have been solved on more conventional CPU clusters before. For example, a massively parallel implementation of a multigrid solver on hybrid grids is described in [34] and the authors demonstrate the excellent scalability of the algorithm on nearly 300,000 CPU cores by solving systems with up to 10^{12} unknowns.

While in the past, bespoke geometric multigrid solvers for anisotropic elliptic PDEs have been studied extensively in the literature (see e.g. [35] for a standard textbook), we will not discuss those more algorithmic aspects here and instead refer the reader to [2] and a forthcoming publication ([15]) which contain more comprehensive reviews of this topic.

Structure. This paper is organised as follows: in Section 2 we briefly review the application of iterative solvers to anisotropic elliptic PDEs in atmospheric modelling with particular focus on Conjugate Gradient and geometric multigrid methods. The GPU implementation of these methods is discussed in Section 3 and a theoretical performance analysis is carried out in Section 4. The results of our numerical experiments and weak scaling tests on Titan are presented in Section 5. Finally we conclude and outline some ideas for further work in Section 6.

2. Iterative solvers for anisotropic elliptic PDEs in implicit time stepping methods

2.1. Model problem

We consider the following PDE, which can be used as a simplified model of the pressure correction equation arising in semi-implicit semi-Lagrangian time stepping methods in atmospheric forecast models:

$$-\omega(h)^2 \left(\nabla_{\mathcal{S}}^2 u(\hat{\mathbf{r}}, r) + \lambda(h)^2 \frac{1}{r^2} \frac{\partial}{\partial r} \left(r^2 \frac{\partial u(\hat{\mathbf{r}}, r)}{\partial r} \right) \right) + u(\hat{\mathbf{r}}, r) = f(\hat{\mathbf{r}}, r). \quad (3)$$

Here $r \in [1, 1 + H]$ is the radial coordinate in units of the earth's radius R_{earth} and $H = D/R_{\text{earth}} \ll 1$ is the ratio between the depth of the atmosphere and the radius of the earth. The unit vector $\hat{\mathbf{r}}$ is used to describe a position on a unit sphere \mathcal{S} and $\nabla_{\mathcal{S}} \equiv \nabla - \langle \hat{\mathbf{r}}, \nabla \rangle \hat{\mathbf{r}}$ denotes the tangential component of the three dimensional gradient. Homogeneous Neumann boundary conditions are used at the top and bottom boundary of the domain. A structured vertical grid and a semi-structured horizontal grid \mathcal{T}_h are used for discretising the equation on the domain $\mathcal{S} \times [1, 1 + H]$. Since $H \ll 1$ the vertical grid spacing h_z is much smaller than the horizontal mesh width h , and so the discretised equation has a very strong grid aligned anisotropy in the vertical direction. The parameters $\omega(h)$ and $\lambda(h)$ depend on the meteorological conditions and on the time step size and are discussed in more details in [1]. Most importantly, as the horizontal resolution increases, i.e. the mesh width h tends to zero, we have $\omega(h) \propto h$ and $\lambda(h) \rightarrow 1$. More specifically the coefficient of the second order term in (3) is given by

$$\omega(h) = \frac{c_h \Delta t(h)}{2R_{\text{earth}}} \quad (4)$$

where c_h is at the order of the speed of sound. Because of fast advective time scales and to represent large scale flow accurately, in meteorological applications the resolution dependent time step size $\Delta t(h)$ has to be chosen such that the

horizontal CFL number $\nu_{CFL} = c_h \Delta t / \Delta x = 2\omega(h)/h$ is not larger than around 2 – 10. On the other hand, implicit time stepping methods will not be competitive if ν_{CFL} is too small. To satisfy these conditions we always use

$$\omega(h) = \frac{1}{2} \nu_{CFL} h \quad \text{with} \quad \nu_{CFL} = 8.4 \quad (5)$$

in our numerical experiments. We also study the robustness of our solvers to variations in the CFL number.

The equation in (3) can be seen as a special case of the more general PDE studied in [15] (see also [36] which describes how an equation of this form is derived in the ENDGame dynamical core of the UK Met Office’s Unified Model)

$$-\omega^2 \nabla \cdot (\underline{\alpha}(\hat{\mathbf{r}}, r) \nabla u(\hat{\mathbf{r}}, r)) - \omega^2 \underline{\xi}(\hat{\mathbf{r}}, r) \cdot \nabla u(\hat{\mathbf{r}}, r) + \beta(\hat{\mathbf{r}}, r) u(\hat{\mathbf{r}}, r) = f(\hat{\mathbf{r}}, r), \quad (6)$$

where $\underline{\alpha}$, $\underline{\xi}$ and β are atmospheric “profiles”, i.e. functions which depend on the current state of the model. Due to the vertical layering of the atmosphere each of these functions can be approximated very well as the product of a vertically varying field and a horizontally varying function. This is why tensor-product methods are of particular interest. Even if the profiles do not factorise exactly, an approximate factorisation can still be used in a preconditioner. This is discussed in a lot of detail in [15]. It is for these reasons that we believe that the PDE in (3) is a good model for the pressure correction equation encountered in atmospheric models.

We chose not to make any further simplifications such as solving a (shifted) Laplace equation in a simplified geometry, as is often done in the literature on massively parallel solvers, since this would allow significant further performance improvements which are not reasonable in realistic meteorological applications.

As discussed below one of the crucial advantages of our matrix-free implementation is that it allows a generalisation of problems in flat Euclidean domains (with constant matrix stencil) to spherical geometries without significant additional computational costs.

2.1.1. Discretisation

Equation (3) is discretised using a simple cell centred finite volume scheme on one panel of a non-conformal cubed sphere grid with gnomonic projection as described in [37].

For simplicity homogeneous Dirichlet boundary conditions are used in the horizontal direction. To represent a field $u(\hat{\mathbf{r}}, r)$, all data in one vertical column above the horizontal grid cell T is stored in a vector $\bar{\mathbf{u}}^{(T)}$ of length n_z . Then the discretised equations associated with the horizontal grid cell $T \in \mathcal{T}_h$ can be written as

$$(\underline{A}\bar{\mathbf{u}})^{(T)} = \underline{A}_T \bar{\mathbf{u}}^{(T)} + \sum_{T' \in \mathcal{N}(T)} \underline{A}_{TT'} \bar{\mathbf{u}}^{(T')} = \bar{\mathbf{f}}^{(T)}. \quad (7)$$

where \underline{A}_T is a tridiagonal matrix containing all vertical couplings, as well as diagonal terms, and where the diagonal matrices $\underline{A}_{TT'}$ describe the couplings to the horizontally neighbouring cells $T' \in \mathcal{N}(T)$. Due to the strong vertical anisotropy, the entries in \underline{A}_T are much larger than the ones in $\underline{A}_{TT'}$.

To understand the origin of the individual terms in (7), it is instructive to write down the explicit form of these matrices for a simplified equation instead of (3). Consider the shifted Laplace equation $-\omega^2 (\nabla_{2D}^2 u + \lambda^2 \partial^2 / \partial_r^2 u) + u = f$ in a flat box $\Omega \times [0, H]$; here the horizontal domain $\Omega = [0, 1] \times [0, 1]$ is the unit square and $\nabla_{2D}^2 = \partial^2 / \partial_x^2 + \partial^2 / \partial_y^2$ denotes the two dimensional Laplacian. If we choose an equidistant Cartesian grid with spacing h on Ω then every horizontal cell T can be labelled with a pair of indices, i.e. $T \equiv (i, j)$. In this case equation (7) can be written explicitly as

$$(\underline{A}\bar{\mathbf{u}})^{(i,j)} = \underline{A}_{(i,j)} \bar{\mathbf{u}}^{(i,j)} + \underline{A}_{(i,j),(i+1,j)} \bar{\mathbf{u}}^{(i+1,j)} + \underline{A}_{(i,j),(i-1,j)} \bar{\mathbf{u}}^{(i-1,j)} + \underline{A}_{(i,j),(i,j+1)} \bar{\mathbf{u}}^{(i,j+1)} + \underline{A}_{(i,j),(i,j-1)} \bar{\mathbf{u}}^{(i,j-1)} = \bar{\mathbf{f}}^{(i,j)}$$

with the matrices $\underline{A}_{(i,j),(i',j')} = -\omega^2 / h^2 \mathbb{I}_{n_z \times n_z}$ where $\mathbb{I}_{n_z \times n_z}$ is the $n_z \times n_z$ identity matrix. The entries on the diagonal of the matrix $\underline{A}_{(i,j)}$ would be $1 + 4\omega^2 / h^2 + 2\omega^2 \lambda^2 / h_z^2$ and the off-diagonal entries are $-\omega^2 \lambda^2 / h_z^2$ where h_z is the vertical grid spacing. We stress, however, that equation (7) allows more general geometries with semi-structured horizontal grids. In this case the exact form of the matrices $\underline{A}_{TT'}$ and \underline{A}_T is more complicated as the finite volume discretisation leads to non-trivial geometric factors.

The elliptic equation in (3) is symmetric and positive definite. For a given horizontal resolution we can give a rough estimate of the condition number by again considering the equation in a flat box. After preconditioning by vertical line relaxation, it can be shown that the resolution dependent condition number $\kappa(h)$ is

$$\kappa(h) \approx 1 + 8\omega(h)^2 / h^2 \quad \text{for} \quad h \ll 1 \quad \text{and} \quad \omega(h) \ll 1. \quad (8)$$

Since $\omega(h) = 4.2h$ for $\nu_{CFL} = 8.4$, this leads to an estimate of $\kappa(h) \approx 142$. Geometric factors arising from the spherical geometry will modify this estimate by factors of $\mathcal{O}(1)$ and hence we expect the condition number of our problem to be in the range 100 – 1000, independent of the grid resolution.

In principle we do not need to make any assumptions on the ordering of the horizontal degrees of freedom and indirect addressing could be used in the horizontal direction, as is described for the DUNE implementation of the problem on quasi-uniform icosahedral and cubed-sphere grids for the entire sphere in [15]. However, in this work we assume for simplicity that each horizontal grid cell on the panel can be identified by a pair of indices $(i, j) \in [1, n_x] \times [1, n_y]$, and each vertical level is indexed by an additional integer $k \in [0, n_z - 1]$.

2.2. Algorithmically scalable and efficient solvers

Large sparse systems of equations can be solved efficiently using state-of-the-art iterative solvers which improve an initial solution \mathbf{u}_0 by reducing the residual $\mathbf{r} = \mathbf{f} - \underline{A}\mathbf{u}$ (and hence the error) at every iteration. Krylov subspace methods (see e.g. [38] for an introduction) minimise the residual by constructing the solution in the space spanned by the vectors

$$\mathbf{r}_0, \underline{A}\mathbf{r}_0, \underline{A}^2\mathbf{r}_0, \dots, \quad (9)$$

where $\mathbf{r}_0 = \mathbf{f} - \underline{A}\mathbf{u}_0$ is the initial residual. The simplest (and most efficient) Krylov subspace method for symmetric positive systems is the preconditioned Conjugate Gradient (CG) iteration. Closely related methods such as Conjugate Residual (CR), GMRES and BiCGStab are very popular in the meteorological literature (see e.g. [8, 9, 10, 11]) and due to the strong vertical anisotropy, a very effective preconditioner \underline{M} is vertical line relaxation, which requires the solution of a tridiagonal problem in each vertical column. This preconditioner corresponds to solving the equation which is obtained by only keeping the first term on the left hand side of (7), which describes the dominant vertical couplings; the resulting matrix \underline{M} is block-diagonal. Each of the tridiagonal systems can be solved independently using the Thomas algorithm (see e.g. [12]). Mathematically this is equivalent to a block-Jacobi or block-SOR method where each of the blocks correspond to the degrees of freedom in one particular vertical column.

The computationally most expensive components of the algorithm are a sparse matrix-vector (SpMV) multiplication and a preconditioner (tridiagonal-) solve, in the following we write these operations as

$$\mathbf{y} \leftarrow \underline{A}\mathbf{x} \quad (\text{SpMV}), \quad \mathbf{y} \leftarrow \underline{M}^{-1}\mathbf{x} \quad (\text{Preconditioner}). \quad (10)$$

As discussed in detail in [13] (where also the algorithm is written down explicitly), the efficiency of the implementation can be improved by fusing these two operations with the level 1 BLAS operations in the main loop. Other Krylov subspace methods, such as BiCGStab, CR or GMRES can be used to solve more general systems and differ from the CG method only in the number of sparse matrix-vector products, preconditioner applications, level 1 BLAS operations, and in the storage requirements.

In contrast, multigrid methods (see e.g. [39, 35]) use a hierarchy of coarse levels to minimise the error on all scales simultaneously. In the following we write $\mathbf{u}^{(\ell)}$ for the field on multigrid level ℓ , where $\ell = 1$ corresponds to the coarsest level and $\ell = L$ to the finest level where we want to solve the equation, i.e. $\mathbf{u}^{(L)} = \mathbf{u}$. For simplicity we omit the multigrid-level index (ℓ) wherever it is obvious from the context, such as on all the coarse grid matrices. The fine grid equation is solved by starting with an initial guess for the solution and improving on this by repeated calls to the recursive subroutine `Vcycle` in Algorithm 1 (for simplicity the iteration is written down for one pre- and one post-smoothing step here). After each `Vcycle` convergence is checked by comparing the norm of the residual to a given tolerance ε , i.e. the algorithm terminates as soon as

$$\|\mathbf{r}\|/\|\mathbf{r}_0\| < \varepsilon. \quad (11)$$

In our numerical experiments we always reduce the residual by five orders of magnitude, which is typical in atmospheric applications. To achieve rapid convergence the different multigrid components have to be adapted to the problem to be solved. In [2] geometric multigrid algorithms for equations with a tensor-product structure and grid-aligned anisotropy are analysed. The authors show that the convergence rate of a multigrid solver for a two-dimensional problem with strong coupling in the vertical direction can be bounded by the convergence rate of the multigrid algorithm for a related one-dimensional horizontal problem if the following tensor-product multigrid algorithm is used:

- **Horizontal-only semicoarsening:** Only coarsen the grid in the horizontal direction.
- **vertical block-Jacobi/-SOR smoother:** Use vertical line relaxation as the smoother, i.e. solve the equation for all degrees of freedom in a vertical column simultaneously. Hence, in essence the multigrid smoother is identical to the preconditioner used for the CG algorithm described above.

As shown in [15], the generalisation from two to three dimensions is straightforward and this is the algorithm which we use in this work. On each level we use a block-Jacobi smoother which can be written as

$$\mathbf{u}^{(\ell)} \leftarrow \mathbf{u}^{(\ell)} + \rho_{\text{relax}}\underline{M}^{-1} \left(\mathbf{f}^{(\ell)} - \underline{A}\mathbf{u}^{(\ell)} \right) \quad (12)$$

Algorithm 1 Subroutine $\text{VCycle}(\rho_{\text{relax}}, \{\mathbf{u}^{(\ell)}\}, \{\mathbf{f}^{(\ell)}\}, \{\mathbf{r}^{(\ell)}\}, \ell)$

```
if  $\ell = 1$  then
  {Restrict residual and solve on coarsest level}
   $\mathbf{f}^{(1)} \leftarrow \underline{R}_{1,2} \mathbf{r}^{(2)}$ ,  $\mathbf{u}^{(1)} \leftarrow \underline{A}^{-1} \mathbf{f}^{(1)}$ 
else
  if  $\ell = L$  then
    {Smooth once on finest level}
     $\mathbf{u}^{(\ell)} \leftarrow \mathbf{u}^{(\ell)} + \rho_{\text{relax}} \underline{M}^{-1}(\mathbf{f}^{(\ell)} - \underline{A}\mathbf{u}^{(\ell)})$  [= Kernel Smooth]
  else
    {On all other levels, restrict residual and smooth once}
     $\mathbf{f}^{(\ell)} \leftarrow \underline{R}_{\ell, \ell+1} \mathbf{r}^{(\ell+1)}$ ,  $\mathbf{u}^{(\ell)} \leftarrow \rho_{\text{relax}} \underline{M}^{-1} \mathbf{f}^{(\ell)}$  [= Kernel RestrictSmooth]
  end if
  {Calculate residual}
   $\mathbf{r}^{(\ell)} \leftarrow \mathbf{f}^{(\ell)} - \underline{A}\mathbf{u}^{(\ell)}$  [= Kernel Residual]
  {Recursive call to Subroutine VCycle}
  Call  $\text{VCycle}(\rho_{\text{relax}}, \{\mathbf{u}^{(\ell)}\}, \{\mathbf{f}^{(\ell)}\}, \{\mathbf{r}^{(\ell)}\}, \ell - 1)$ 
  {Add prologated coarse grid correction}
   $\mathbf{u}^{(\ell)} \leftarrow \mathbf{u}^{(\ell)} + \underline{P}_{\ell, \ell-1} \mathbf{u}^{(\ell-1)}$  [= Kernel Prolongate]
  {Postsmoothing}
   $\mathbf{u}^{(\ell)} \leftarrow \mathbf{u}^{(\ell)} + \rho_{\text{relax}} \underline{M}^{-1}(\mathbf{f}^{(\ell)} - \underline{A}\mathbf{u}^{(\ell)})$  [= Kernel Smooth]
end if
```

and requires one sparse-matrix-vector product and one preconditioner solve in (10). For the intergrid operations

$$\mathbf{u}^{(\ell)} \leftarrow \underline{R}_{\ell, \ell+1} \mathbf{u}^{(\ell+1)} \quad (\text{Restriction}), \quad \mathbf{u}^{(\ell)} \leftarrow \underline{P}_{\ell, \ell-1} \mathbf{u}^{(\ell-1)} \quad (\text{Prolongation}) \quad (13)$$

we use a simple cell-average for the restriction and (piecewise) linear interpolation for prologation (both in the horizontal direction only), and we found that these methods are sufficient for scalable performance. By carrying out the restriction at the beginning of the subroutine in Algorithm 1 it is possible to fuse it with the first presmoothing step on the coarse levels. Apart from that fusing kernels has little potential for further gains in the multigrid algorithm.

Recall that the condition number κ_{fine} of the fine grid problem is $\mathcal{O}(100-1000)$ independent of the horizontal resolution. The condition number of each subsequent coarse level is reduced by a factor 4, i.e. the square of the relative grid spacings. We typically choose $L = 5$ multigrid levels. Hence on the coarsest grid we have $\kappa_{\text{coarse}} = 4^{-(L-1)} \kappa_{\text{fine}}$, i.e. the operator is well conditioned and the coarse grid equation can be solved by a small number (two turned out to be sufficient) of smoother iterations. This has already been confirmed by the detailed numerical experiments on CPU architectures in [1]. For the particular model problem in (3) where $\omega(h) \propto h$ to accurately represent large scale atmospheric flow, both these iterative methods are algorithmically scalable, i.e. the number of iterations is independent of the mesh size h and thus of the problem size. However, an additional benefit of multigrid solvers is their greater robustness with respect to variations in the model coefficients [1].

3. Implementation

In the following we describe the CUDA-C implementations of both the CG- and multigrid solvers discussed in the previous section. The code was written from scratch by the authors and we use the CUBLAS library for some level 1 BLAS operations as well as the GCL library [14] for inter-GPU communication. The source code is made available under the LGPL 3 license and can be accessed as a git repository via the following link: <https://bitbucket.org/em459/ellipticsolvergpu>. Further details on the single GPU implementation of our CG solver can be found in [13].

3.1. Memory throughput optimised implementation

Due to the vertical dependency in the tridiagonal solver, which is used both as the preconditioner in CG and as the smoother in multigrid, one thread is assigned to each vertical column. To achieve optimal performance it is crucial to coalesce access to global memory for all threads within one warp. This is achieved by storing data contiguously in the horizontal (x -) direction. As discussed in Sections 3.2 and 6.1 of [13], a three dimensional field \mathbf{u} on one panel of a cubed

sphere grid can be described as a collection of $n_x \times n_y$ vectors $\bar{\mathbf{u}}^{(i,j)}$, one for each horizontal cell $(i, j) \in [1, n_x] \times [1, n_y]$. Internally the field can be stored as a linear array of length $n_x \times n_y \times n_z$ defined by the mapping,

$$u_{\Lambda(i,j,k)} = \bar{u}_k^{(i,j)}, \quad \text{where} \quad \Lambda(i, j, k) \equiv n_x \cdot (n_z \cdot (j - 1) + k) + (i - 1). \quad (14)$$

In the following we also assume that (at least on the finest multigrid levels) both n_x and n_y are multiples of 32. This further improves performance as global memory access is not only coalesced but also well-aligned. We stress, however, that our approach can be generalised to more unstructured grids e.g. by the use of a space filling curve for numbering the horizontal grid cells.

Matrix-free implementation. As in [13] we use a matrix-free implementation, i.e. we recalculate the local matrix stencil whenever it is needed. In particular, the diagonal matrix $\underline{A}_{T,T'}$ and the tridiagonal matrix \underline{A}_T in (7) are given by

$$\begin{aligned} \underline{A}_{T,T'} &= \alpha_{T,T'} \text{diag}(\mathbf{d}), \\ \underline{A}_T &= |T| \text{diag}(\mathbf{a}) - \alpha_T \text{diag}(\mathbf{d}) + |T| \text{tridiag}(-(\mathbf{b} + \mathbf{c}), \mathbf{b}, \mathbf{c}). \end{aligned} \quad (15)$$

The four vectors \mathbf{a} , \mathbf{b} , \mathbf{c} and \mathbf{d} have length n_z and can be derived from the vertical stiffness- and mass- matrices in a tensor product representation of (3) [15]. As they do not differ from column to column they can be precomputed once for the entire grid. The coefficients $\alpha_{T,T'}$ and α_T are different for each horizontal grid cell T (and depend on the multigrid level). However, they are scalars which can be computed for an entire vertical column with a small overhead as long as n_z is sufficiently large (in atmospheric applications $n_z = \mathcal{O}(100)$, and we use $n_z = 128$ throughout this work). This should be compared to a matrix-explicit implementation where seven matrix entries need to be loaded from memory per grid cell to carry out a sparse matrix-vector product. The matrix-free implementation significantly reduces global memory access, in particular if the vectors \mathbf{a} , \mathbf{b} , \mathbf{c} and \mathbf{d} are cached. Instead of $\mathcal{O}(7 \times n_{\text{horiz}} \times n_z)$ the storage requirements of the matrix are only $\mathcal{O}(4 \times n_z)$ which also means that significantly larger problems can be solved on a single GPU.

At this point we mention that a very similar structure for the matrices $\underline{A}_{T,T'}$ and \underline{A}_T arises when discretising more general equations of the form

$$-\omega^2 \nabla \cdot (\underline{\alpha} \nabla u) - \omega^2 \boldsymbol{\xi} \cdot \nabla u + \beta u = f$$

as long as the ‘‘profiles’’ $\underline{\alpha}$, β and $\boldsymbol{\xi}$ can be factorised into a parts which contain on the horizontal and vertical coordinates only. In this case the factors $|T|$, $\alpha_{T,T'}$ and α_T in (15) have to be replaced by more complicated expressions which can, however, still be evaluated only once per column. It is for this reason that we believe that the results in this article can be generalised easily to more complicated equations, this is discussed in more detail in [15].

On the other hand, in a flat, Euclidean geometry the matrix stencil would be the same for all grid cells and in the simplest possible implementation (which just uses a constant matrix stencil) there would be no cost associated with the local matrix assembly. However, as argued above, for large enough n_z our approach also does not require any additional computational cost (measured in the amount of data read from memory) for setting up the local matrix in each cell. Therefore it allows the generalisation of the problem to spherical domains at the same computational cost as is required for the simplest possible implementation in a flat, Euclidean domain.

In addition, in the CG algorithm we reduce the amount of global memory access by fusing the two computationally most expensive kernels (SpMV and tridiagonal solve) with several of the level 1 BLAS operations. The following operations need to be carried out for the solvers and were implemented as CUDA-C kernels:

Conjugate Gradient.

1. *Sparse matrix-vector product* [Kernel (Fused) SpMV]: Simultaneously calculate $\mathbf{u} \leftarrow \mathbf{u} + \alpha \mathbf{p}$; $\mathbf{p} \leftarrow \mathbf{z} + \beta \mathbf{p}$; $\mathbf{q} \leftarrow \underline{A} \mathbf{z} + \beta \mathbf{q}$; $\sigma \leftarrow \langle \mathbf{p}, \mathbf{q} \rangle$
2. *Preconditioner (tridiagonal solve)* [Kernel (Fused) Tridiag]: Simultaneously calculate $\mathbf{r} \leftarrow \mathbf{r} - \alpha \mathbf{q}$; $\mathbf{z} \leftarrow \underline{M}^{-1} \mathbf{r}$; $\|\mathbf{r}\| \leftarrow \langle \mathbf{r}, \mathbf{r} \rangle$; $\kappa \leftarrow \langle \mathbf{r}, \mathbf{z} \rangle$

Multigrid.

1. *Smoother* [Kernel Smooth]: $\mathbf{u}^{(\ell)} \leftarrow \mathbf{u}^{(\ell)} + \rho_{\text{relax}} \underline{M}^{-1} (\mathbf{f}^{(\ell)} - \underline{A} \mathbf{u}^{(\ell)})$; to avoid a race condition in $\mathbf{u}^{(\ell)}$, this operation is split up into two kernels, with the residual calculation and forward sweep of the tridiagonal solver in the first and the backward sweep and axpy-like update of \mathbf{u} in the second.
2. *Residual calculation* [Kernel Residual]: $\mathbf{r}^{(\ell)} \leftarrow \mathbf{f}^{(\ell)} - \underline{A} \mathbf{u}^{(\ell)}$
3. *Interleaved (fused) restriction and smoother* [Kernel RestrictSmooth]: Simultaneously calculate $\mathbf{f}^{(\ell)} \leftarrow \underline{R}_{\ell, \ell+1} \mathbf{r}^{(\ell+1)}$, $\mathbf{u}^{(\ell)} \leftarrow \rho_{\text{relax}} \underline{M}^{-1} \mathbf{f}^{(\ell)}$

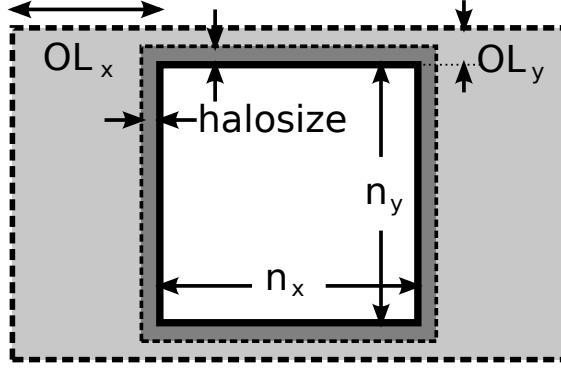


Figure 1: Data layout of the local subdomain on one processor. Interior degrees of freedom $(i, j) \in [1, n_x] \times [1, n_y]$ “owned” by the processor are shown in white, the halo is shown in dark gray and extra padding space in light gray. We set $OL_x = 32$ and $OL_y = \text{halosize}$ to guarantee aligned memory access in the x -direction.

4. Prologation [Kernel Prolongate]: $\mathbf{u}^{(\ell)} = \mathbf{u}^{(\ell)} + \underline{P}_{\ell, \ell-1} \mathbf{u}^{(\ell-1)}$

Each of the kernels requires a single iteration over the grid. For the multigrid solver the first presmoothing step on the coarse grids (on which the initial solution is initialised to zero) has been fused with the residual calculation to reduce the number of accesses to the vector $\mathbf{f}^{(\ell)}$.

In addition, a small number of level 1 BLAS operations still needs to be carried out for global reductions. For example, the global reductions in the interleaved CG kernels are implemented by each thread summing up the values in a vertical column into a two-dimensional field, which is then summed by a cuBLAS dot product with a field that is set to 1 in the whole two dimensional domain. For simplicity the calculation of the norm of the residual on the finest multigrid level was also implemented via a 3D cuBLAS norm instead of fusing it with the corresponding kernel (note that this norm calculation is also not necessary if a fixed number of V-cycles is carried out). While there is further potential for (small) additional speedups, we found that the cost of these level 1 BLAS operations is negligible (to see this, compare the last two rows in Tables 4 and 5 below) and not worth the effort.

3.2. Multi-GPU implementation

To parallelise the solvers across several GPUs we split the horizontal domain into equal square parts, such that each GPU is responsible for a quadratic subdomain. Consistency between neighbouring domains is guaranteed by exchanging halo data when necessary. For this we used the Generic Communication Library (GCL) [14]. The extension of this approach to more general partitionings and different (semi-) structured horizontal grids is possible. This is discussed in [15] where we discuss a parallel CPU implementation of the tensor-product multigrid solver on an icosahedral grid.

In addition to the interior degrees of freedom with horizontal indices $(i, j) \in [1, n_x] \times [1, n_y]$ a halo of cells of width $\text{halosize} = 1$ is stored on each GPU. To avoid unaligned memory access, which we found can reduce performance by as much as 30%, we pad data by a total amount of $OL_x = 32$ in the x -direction, and set $OL_y = \text{halosize}$ in the y -direction (in this direction padding is not necessary). Then the linear mapping (14) is modified to

$$\hat{\mathbf{A}}(i, j, k) \equiv (n_x + 2OL_x) \cdot (n_z \cdot (j - 1 - OL_y)) + (i - 1 - OL_x) \quad (16)$$

with $(i, j) \in [1 - OL_x, n_x + OL_x - 1] \times [1 - OL_y, n_y + OL_y - 1]$; the local domain is shown in Figure 1. We stress, however, that only degrees of freedom on the halo cells are exchanged between processors i.e. the padding does not increase the amount of data that is sent over the network. In the GCL this can be achieved by registering the appropriate degrees of freedom in the halo exchanger object

```

he->add_halo<0>(halosize, halosize, OL_X, NX+OL_X-1, NX+2*OL_X);
he->add_halo<1>(halosize, halosize, OL_Y, NY+OL_Y-1, NY+2*OL_Y);
he->add_halo<2>(0, 0, 0, NZ-1, NZ);

```

where `he` is an instantiation of the *uniform type* halo structure interface class `GCL::halo_exchange_dynamic_ut`.

(Fused) CG				Multigrid			
	FLOPs	Mem	Mem ^(C)		FLOPs	Mem	Mem ^(C)
(Fused) SpMV	32	12	6	Smooth	37	17	8
(Fused) Tridiag	22	12	9	RestrictSmooth	17	12	6
Total	54	24	15	Residual	23	9	3
				Prolongate	6	5	3
				Residual norm	2	1	1
				Total	149.4	67.2	29.6
				(fine level only)	(122)	(53)	(23)

Table 1: Number of FLOPs and memory references per grid cell for the kernels in the (fused) CG algorithm (left) and the multigrid solver (right); see Section 3.1 and Algorithm 1 for a definition of the individual kernels.

4. Theoretical performance analysis

4.1. Floating point operations and memory transfer costs

The number of floating point operations (FLOPs) and memory references per grid cell is shown for the CG and multigrid algorithms in Table 1. The total number of operations on all 5 multigrid levels (penultimate row in Table 1, right) was obtained by adding up the number of operations on all levels and dividing by the size of the fine grid. We assume that 1 pre- and post-smoothing step is applied on each level and the coarse grid problem is solved by two smoother iterations. The residual norm is only calculated on the fine level to test for convergence.

The column labelled with “Mem” shows the number of global memory references without any caching. On the other hand the “Mem^(C)” column shows the corresponding number assuming optimal caching, i.e. we assume that data is only loaded from global memory once per kernel call. As our implementation is matrix-free (and we assume that the vectors \mathbf{a} , \mathbf{b} , \mathbf{c} and \mathbf{d} in (15) are always cached), there are no costs associated with reading the local stencil from global memory. Consider, for example, the sparse matrix vector product $\mathbf{y} \leftarrow \underline{A}\mathbf{x}$. Without caching, the value of the vector \mathbf{x} in each cell and its six direct neighbours needs to be read from memory and one value is written back to \mathbf{y} , resulting in a total of 8 memory references. With perfect caching each entry in \mathbf{x} only has to be read from global memory once and the total number of memory references is reduced to 2 per grid cell. How much caching can actually be achieved is hard to predict, so these two values should be interpreted as upper and lower bounds. As explained in Section 5.1 we always use the lower value (i.e. the number in the “Mem^(C)” column) for our estimates of the achieved global memory bandwidth. The quantity which is constructed in this way is commonly known as the “useful bandwidth” since it does not include any spurious memory traffic which is not required by the algorithm (for example data which is read twice due to poor caching).

In the multigrid solver the grid size is reduced by a factor 4 with each coarsening step and as a result most of the computational work is concentrated on the finest grid level. To demonstrate this we also list the number of operations per grid cell on the fine level only in the last row of Table 1, right. In practise, once the size of the horizontal grid drops below a certain threshold, the GPU might not be utilised efficiently and the actual runtime is reduced by less than a factor four on subsequent levels. Our measurements (see Table 3) show that calculations on the fine grid account for most of the total computational cost and that the cost reduction factor is close to four as long as the local domain is larger than 64×64 (see Figure 3 (left)). The impact of inter-GPU communications, which is also more significant on the coarser grids, will be discussed below.

For the CG solver, the number of FLOPs is only $3.6\times$ larger than the minimal number of memory references, for the multigrid solver the corresponding factor is $5\times$, and we conclude that both algorithms are clearly memory bound on a GPU. By comparing the number of memory references in the two algorithms, the theoretically expected time per iteration is $2\times$ to $3\times$ larger for the multigrid solver where the exact ratio depends on the cache efficiency.

It is worth counting the additional number of global memory accesses which would be required in a matrix-explicit code: to read the matrix from memory requires 10 memory accesses for the CG solver and 42.9 for the multigrid algorithm (the full seven point stencil is required to evaluate the sparse matrix-vector product, but only the couplings to the cell above and below are required in the tridiagonal solver). This should be compared to the minimal number of memory references shown for the matrix-free code in Table 1 which is 15 for CG and 29.6 for multigrid.

4.2. Parallel communications between GPUs

In the CG solver a halo exchange is required after each call to the (fused) tridiagonal solver kernel. In the multigrid algorithm halos need to be exchanged after each kernel launch (with the exception of the residual calculation). Denoting

Solver	Problem size n_x			
	128	256	512	768
(Fused) CG	20.8	10.4	5.2	3.5
Multigrid	80.5	40.2	20.1	13.4
Multigrid (fine level)	54.3	27.2	13.6	9.1

Table 2: Theoretical ratio between MPI communication- and calculation- times as defined in (17). The table shows R multiplied by 10^4 for different local horizontal problem sizes $n_x = n_y$ with `halosize` = 1 and double precision arithmetic (`sizeof(double)` = 8).

the number of halo exchanges per iteration by n_{halo} and the minimal number of memory references by $n_{\text{Mem}(C)}$, on a given level with a local problem size of $n_x \times n_y \times n_z$ (where we always implicitly assume that $n_x = n_y$), the ratio between the communication time and calculation time is given by

$$\rho = \frac{2(n_x + n_y)n_z \times \text{sizeof}(\text{double}) \times \text{halosize} \times n_{\text{halo}}}{n_x n_y n_z \times n_{\text{Mem}(C)}} \times \frac{BW_{\text{mem}}}{BW_{\text{MPI}}} \equiv R \times \frac{BW_{\text{mem}}}{BW_{\text{MPI}}}. \quad (17)$$

This ratio decreases as $\propto (n_x + n_y)/(n_x \cdot n_y) \propto 1/n_x$ as the local domain size n_x increases. We assume that the global memory bandwidth BW_{mem} is about two orders of magnitude larger than the network bandwidth BW_{MPI} required for communication between different GPUs (the exact ratio between the bandwidths will be quantified in more detail in Section 5.3). For the multigrid solver the amount of exchanged data and the number of memory references have to be summed over all levels to calculate the ratio R in (17).

In Table 2 the ratio R is shown for all solvers. For the multigrid algorithm we also distinguish between the full level hierarchy and the finest level only. While (based on our estimate $BW_{\text{mem}}/BW_{\text{MPI}} \approx 100$) the resulting ratio ρ is at the order of 10% or less for the CG solver and on the finest multigrid level (for local problem sizes larger than $256 \times 256 \times 128$), it is approximately four times larger for the multigrid solver. The main reason for this is that three halo exchanges are required on each level instead of one for the CG solver. In addition on the coarser levels the ratio between communications and calculations is a factor $n_x/n_x^{(\ell)} > 1$ larger than on the finest level, where $n_x^{(\ell)} = n_y^{(\ell)}$ is the problem size on level ℓ . With $L = 5$ multigrid levels this corresponds to an increase by a factor $2^{L-1} = 16$ on the coarsest level. For both reasons we expect the overhead from MPI communications to be larger for the multigrid solver and this is confirmed by our numerical experiments in Section 5.3.

5. Results

In the following we report on detailed numerical experiments with our solvers and quantify both their algorithmic- and computational performance on up to 16384 nVidia GPUs.

Hardware and compilers. Most GPU results shown in this section were obtained on the Cray XK7 Titan supercomputer at Oak Ridge National Lab. According to the June 2014 release of the `top500.org` list [6] this machine is currently ranked as the second fastest computer in the world and consists of 18,688 nodes in total. Each node contains one 16-core 2.2GHz AMD Opteron(TM) 6274 (Interlagos) CPU and one nVidia Tesla K20X card with a GK110 Kepler GPU (compute capability 3.5). The nodes are linked with a Cray Gemini interconnect in a torus topology. Each GPU has 2688 cores, which are organised into 15 streaming multiprocessors of 192 cores each. The K20X card has 6 GB of global memory, in addition to 1536KB L2 cache and 48KB+16KB configurable L1 cache/shared memory. The theoretical peak FLOP rate of a single GPU is 1.31 TFLOPs and the theoretical peak global memory bandwidth is 250 GByte/s [7], resulting in a theoretical peak floating point performance of 27 PFLOPs (combined CPU and GPU performance); the LinPACK benchmark performance is quoted as 17.59 PFLOPs [6]. The code was compiled with release 5.0 (V0.2.1221) of the nVidia `nvcc` compiler and version 4.7.2 of the `gnu c++` compiler; a vendor optimised MPICH2 implementation that supports GPUDirect was used for internode communications.

To study the dependency of the performance on the GPU hardware we also carried out a number of runs on up to 64 GPUs of the EMERALD GPU cluster hosted at Rutherford Appleton Lab in the UK. This cluster consists of 372 NVidia M2090 cards with Fermi GPUs (compute capability 3.0), which are organised into 60 nodes with 3 GPUs each and 24 nodes with 8 GPUs each. With 0.67TFLOPs the double precision floating point performance is half as large as for the Kepler GK110 cards in Titan. The global memory bandwidth is 177 GByte/s. Each node contains two 6-core Intel X5650 Xeon CPUs and a QDR Infiniband network is used to link the nodes; the MVAPICH2 implementation was used for MPI communications.

To quantify the performance of our algorithms on traditional CPU hardware we compared the GPU code to a bespoke Fortran implementation of the multigrid solver which is based on the MPI distributed memory programming model. The Fortran code was run on up to 65536 cores of HECToR, the UK’s national supercomputer which is hosted and managed by the Edinburgh Parallel Computing Centre (EPCC). The performance and scalability of this Fortran code is discussed in more detail in [1]. HECToR is a Cray XE6 supercomputer consisting of 2816 compute nodes connected by a Cray Gemini interconnect. Each node contains two AMD Opteron (Interlagos, model 6276, 2.3 GHz) processors with 16 cores [40], and we adjusted the problem size such that it is the same on one GPU and on one CPU, i.e. we always carried out a socket-to-socket comparison.

Problem- and solver- parameters. For a given problem size, which in the following is defined by the number of horizontal grid cells n_x in one direction, the parameters $\omega(h)$ and $\lambda(h)$ in (3) were adjusted to physically realistic values, for the exact values compare Table 6 to Table 1 in [1]. In particular $\omega(h)$ decreases linearly with increasing model resolution as described in (5), which implies that the condition number of the discretised equations does not increase with problem size, and hence the number of iterations is roughly independent of n_x . As estimated above, the condition number is $\mathcal{O}(100 - 1000)$. The number of vertical grid levels was kept fixed at $n_z = 128$.

We use five multigrid levels in all cases and solve the coarse grid problem by applying two smoother iterations and use one pre- and post-smoothing step on all multigrid levels. The overrelaxation parameter in the block-Jacobi smoother was set to $\rho_{\text{relax}} = 2/3$.

5.1. Single GPU performance

The performance of CG and of the multigrid solver is shown in Table 3 for different problem sizes n_x . The largest problem that could be solved with the CG algorithm has $768 \times 768 \times 128 = 7.55 \cdot 10^7$ degrees of freedom, whereas for the multigrid solver the largest problem that fits into memory has $512 \times 512 \times 128 = 3.36 \cdot 10^7$ unknowns. On the GPU the horizontal thread layout was chosen such that each block consists of 64 threads in the (memory contiguous) x -direction and 2 threads in the y -direction. We iterate in both cases until the residual has been reduced by five orders of magnitude, i.e. we use $\varepsilon = 10^{-5}$ in (11). The number of iterations is almost 8 times smaller for the multigrid solves. The numbers depend only weakly on the problem size, confirming the good algorithmic scalability of both solvers. For the $n_x = 512$ problem, one iteration of the multigrid solver is three times as expensive as a CG iteration which is in line with the prediction based on the number of memory references in Section 4.1. This ratio deteriorates for smaller problem sizes because as the problem size decreases, the GPU will be underutilised on the coarse grid levels, see also Figure 3.

At this point it is worth recalling the significant reduction in runtime which is achieved by not storing the matrix explicitly. In [13] the performance of the matrix-free fused CG solver is compared to an implementation based on the CUSparse library which stores the matrix in compressed sparse row (CSR) storage format. For a $256 \times 256 \times 128$ problem on a single Fermi M2090 GPU a speedup of a factor $4.6\times$ is achieved by recalculating the local matrix stencil on-the-fly instead of using the CSR implementation.

In Table 3 we also quantify the floating point performance and the percentage of the peak global memory bandwidth which can be utilised by our solvers. For this we use the figures in Table 1 and divide them by the measured time per iteration in the first row of Table 3. In all cases our calculation is based on the minimal number of memory accesses, i.e. the number in the “Mem^(C)” column, i.e. we calculate the “useful bandwidth”. For a memory bound application a bandwidth of 100% would correspond to perfect caching. It is very hard to achieve this theoretical upper bound in practice.

We stress that the total solution time in the last row includes the time for copying the right hand side from the host to the device and for copying the final solution back to the host. This time is listed separately in Table 3 together with the time for transposing the fields from a z -contiguous data format on the host to the x -contiguous format in (16) on the GPU. For optimal efficiency the transposition was implemented on the GPU by adapting the algorithm described by Mark Harris [41]. Looking at the 512 problem again, the multigrid solver converges twice as fast as the CG iteration. This speedup is reduced for smaller problem sizes but still $1.6\times$ for the smallest problem we considered. In summary this stresses the point made at the beginning of this article: the highest performance gains can be achieved by choosing the algorithmically most efficient solver even if it is more expensive and computationally less efficient in a single iteration and has a slightly worse parallel efficiency.

5.2. Robustness

For Krylov solvers the algorithmic performance, i.e. the number of iterations to reduce the relative residual below a certain threshold, depends on the condition number $\kappa(h)$ of the preconditioned elliptic operator. As discussed above, in the case of vertical line relaxation, $\kappa(h)$ only depends on the ratio of $\omega(h)$ and the grid spacing h ; the size of $\kappa(h)$ is estimated in (8). The dimensionless quantity $\omega(h)$ is proportional to the time step size (see (4)) and hence $2\omega(h)/h = \nu_{CFL}$ is the

	(Fused) CG				Multigrid		
	Problem size n_x				Problem size n_x		
	128	256	512	768	128	256	512
t_{iter}	2.8	7.5	29.7	65.0	12.8	27.3	88.8
GFLOPs	41.0	60.1	60.9	62.7	24.6	45.9	56.4
percentage of peak BW	36.4%	53.4%	54.2%	55.7%	15.6%	29.1%	35.8%
# iterations	70	62	59	58	9	8	8
$t_{\text{MemCpy+transpose}}$	19.4	63.5	228.4	494.0	19.6	64.3	229.5
total solution time	217.1	543.7	2029.0	4365.0	135.5	285.7	949.9

Table 3: Time per iteration t_{iter} , number of iterations and total solution time for different problem sizes n_x and solvers. The total solution time includes the host-device memory transfer and data transposition time, which is listed separately as $t_{\text{MemCpy+transpose}}$. All times are given in milliseconds. The floating point performance and percentage of theoretical peak global memory bandwidth are calculated based on t_{iter} and the numbers in Table 1, see main text for details.

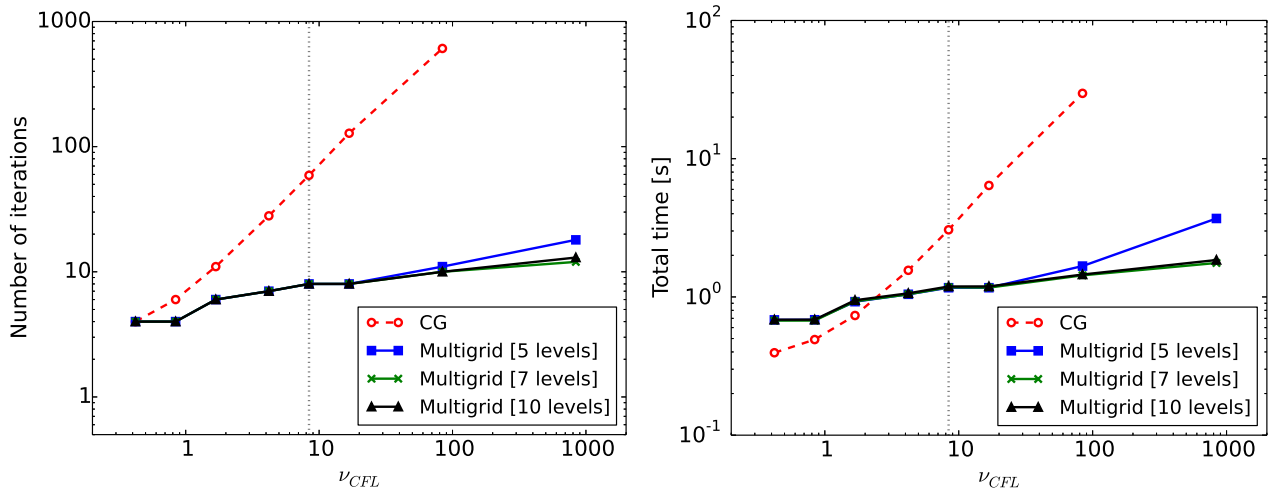


Figure 2: Number of iterations (left) and total solution time (right) for different solvers and a range of CFL numbers. The problem size was fixed to $512 \times 512 \times 128$ and all runs were carried out on one M2090 GPU of the EMERALD cluster. The dashed vertical line marks the CFL number used for numerical experiments in the rest of this article.

horizontal CFL number. As already shown in [1], the tensor-product multigrid method converges in a constant number of iterations for our model problem, independent of the condition number. While for most numerical experiments presented in this paper we fixed $\nu_{\text{CFL}} = 8.4$, here we also study the algorithmic performance for other choices of ν_{CFL} for fixed horizontal resolution. In the limit $\Delta t \rightarrow \infty$, $\nu_{\text{CFL}} \rightarrow \infty$ one would recover the Laplace equation from (3). In Fig. 2 both the number of iterations and the total solution time are shown for different values of ν_{CFL} . In addition to the results for 5 multigrid levels, which were used for the numerical experiments in the rest of this paper, we also show results for 7 and 10 multigrid levels. As expected, since $\kappa(h) \propto \nu_{\text{CFL}}^2$ for $\nu_{\text{CFL}} \gg 1$, the number of CG iterations increases linearly with ν_{CFL} for the CG solver. In terms of the total solution time CG only becomes competitive for small CFL numbers. The multigrid algorithm is significantly more robust. The number of iterations and the total solution time depends only weakly on the CFL number. As the CFL number increases, the coarse grid problem becomes less well conditioned. Depending on the number of levels, it may be necessary to increase the number of smoother iterations on the coarsest grid (or use a different coarse grid solver, such as preconditioned CG). We found that for the 10-level method 2 smoother iterations are sufficient to solve the coarse grid problem up to a value of $\nu_{\text{CFL}} = 840$ without any significant increase in the number of iterations or in the total solution time. For the 7-level method it was necessary to increase the number of smoother iterations on the coarsest level slightly for $\nu_{\text{CFL}} > 16.8$. To maintain robustness, 5 and 15 smoother steps were sufficient for $\nu_{\text{CFL}} = 84$ and $\nu_{\text{CFL}} = 840$, respectively. For the 5-level method, still 2 smoother steps on the coarsest level suffice up to $\nu_{\text{CFL}} = 16.8$, but the number has to increase significantly faster for larger CFL numbers, namely to 30, for $\nu_{\text{CFL}} = 84$, and to 150, for $\nu_{\text{CFL}} = 840$, leading to a slightly faster increase of the number of iterations and of the total solution time

	1 GPU	64 GPUs
(Fused) SpMV	14.1	14.2
(Fused) Tridiag	15.6	18.0
Total [kernels]	29.7	32.2
Total [iteration]	29.8	32.5

Table 4: Breakdown of the (fused) CG solver time per iteration on a $512 \times 512 \times 128$ grid for 1 and 64 GPUs [all times in milliseconds].

kernel	1 GPU					64 GPUs				
	Multigrid level					Multigrid level				
	5	4	3	2	1	5	4	3	2	1
Smooth	45.0	5.9	1.9	0.8	0.7	50.7	7.6	3.2	1.8	1.7
ResSmooth	—	4.7	1.7	0.7	0.7	—	6.3	3.0	1.8	1.6
Residual	15.4	2.0	0.8	0.3	—	21.7	3.5	2.0	1.3	—
Prolongate	4.3	1.1	0.4	0.2	—	6.7	2.7	1.6	1.3	—
Total [kernels]	64.7	13.7	4.8	2.0	1.4	79.1	20.1	9.8	6.3	3.3
Total [iteration]	88.8					122.0				

Table 5: Breakdown of the multigrid solver time per iteration on a $512 \times 512 \times 128$ grid for 1 and 64 GPUs [all times in milliseconds].

for the 5-level method. Nevertheless, with these modifications the multigrid solver is robust over a very wide range of CFL numbers.

5.3. Communication overhead between GPUs and multigrid performance

To identify the bottlenecks of the multigrid algorithm and also to study the impact of parallel communications, the time per iteration was broken down into the time spent in the individual kernels both for a 1 GPU run and for a 64 GPU run with identical local grid size. For the (fused) SpMV and preconditioner kernels these times are shown in Table 4 for a local $512 \times 512 \times 128$ grid. The total time per iteration increases by around 10% due to the halo exchange, this should be compared to the theoretically predicted increase of around 5% according to Table 2. The corresponding results are shown for the multigrid solver in Table 5. Here the time per iteration grows by nearly 40% when going from one to 64 GPUs, while the theoretical analysis predicts a 20% increase.

For the theoretical estimates of the communication/calculation ratio R in (17) we quantify the memory and inter-GPU bandwidth as follows: the peak global memory bandwidth of the K20X card is $BW_{\text{mem}} = 250\text{GB}/s$, and as demonstrated in Section 5.1, our solvers can typically utilise at the order of 30% – 50% of this peak value on a single GPU (see Table 3). We measured the communication bandwidth by carrying out 1000 halo exchanges on 64 GPUs and obtained $BW_{\text{MPI}} \approx 1\text{GB}/s$, which implies that $BW_{\text{mem}}/BW_{\text{MPI}} = \mathcal{O}(100)$.

Looking at the time spent on different multigrid levels, which is also plotted on a logarithmic scale in Figure 3, we see that part of this poor scaling can be attributed to a worse calculation/communication ratio on the coarser levels. On a single GPU the times decrease by roughly a factor of 4 from level to level as expected, until a horizontal problem size of 64×64 is reached. Beyond this point the costs do not decrease further as the GPU is underutilised. However, on 64 GPUs the cost per level is reduced by less than a factor 4 on all levels due to the worse communication and calculation ratio.

5.4. Massively parallel scaling on GPU and CPU clusters

We finally carried out a series of weak scaling runs on Titan. For each of these runs the local problem size was kept fixed at $n_x \times n_y \times n_z$, such that the global problem size grows linearly with the number of GPUs. The resulting numbers of cores and global problem sizes are shown in Table 6. The largest problem solved has half a trillion ($0.55 \cdot 10^{12}$) unknowns. The weak scaling of the time per iteration and the total solution time is plotted in Figure 4. After a slight initial increase all GPU solvers scale very well up to 16384 GPUs and, as expected, the scalability increases with the problem size due to the resulting better calculation / communication ratio. In Figure 5 we compare the performance of the multigrid

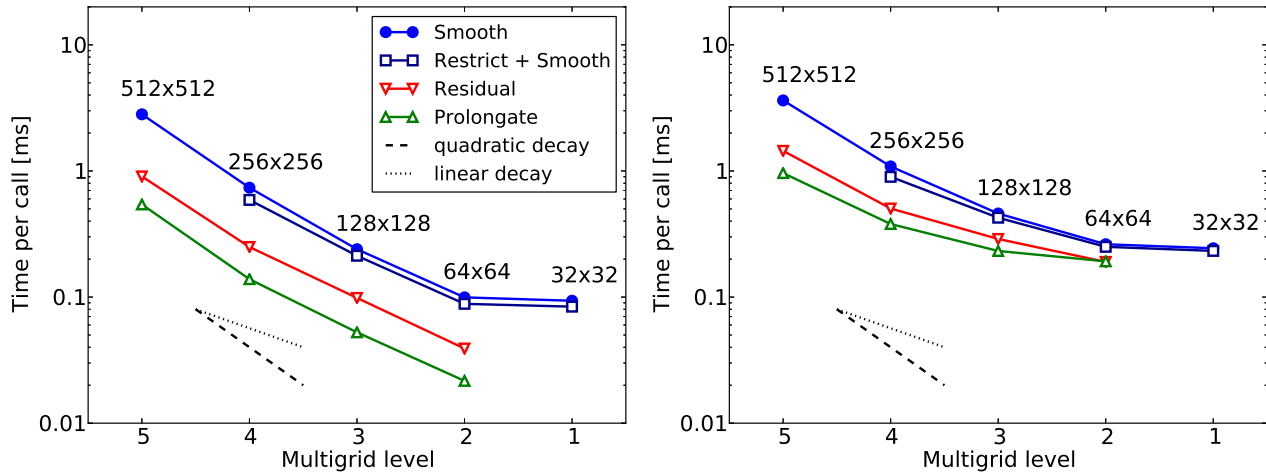


Figure 3: Breakdown of the multigrid solver time on a $512 \times 512 \times 128$ grid on 1 (left) and 64 GPUs (right). The horizontal grid size is shown separately for each multigrid level.

# sockets (p)	#GPU cores		# CPU cores	Local problem size n_x		
	Titan	EMERALD	HECToR	256	512	768
1	2,688	512	16	$8.39 \cdot 10^6$	$3.36 \cdot 10^7$	$7.55 \cdot 10^7$
4	10,752	2,048	64	$3.36 \cdot 10^7$	$1.34 \cdot 10^8$	$3.02 \cdot 10^8$
16	43,008	8,192	256	$1.34 \cdot 10^8$	$5.37 \cdot 10^8$	$1.21 \cdot 10^9$
64	172,032	32,768	1,024	$5.37 \cdot 10^8$	$2.15 \cdot 10^9$	$4.83 \cdot 10^9$
256	688,128	—	4,096	$2.15 \cdot 10^9$	$8.59 \cdot 10^9$	$1.93 \cdot 10^{10}$
1024	2,752,512	—	16,384	$8.59 \cdot 10^9$	$3.44 \cdot 10^{10}$	$7.73 \cdot 10^{10}$
4096	11,010,048	—	65,536	$3.44 \cdot 10^{10}$	$1.37 \cdot 10^{11}$	$3.09 \cdot 10^{11}$
16384	44,040,192	—	—	$1.37 \cdot 10^{11}$	$5.50 \cdot 10^{11}$	—

Table 6: Global number of degrees of freedom $p \times n_x \times n_y \times n_z$ for different numbers of processors p and local problem sizes $n_x = n_y$. The number of vertical columns is always $n_z = 128$.

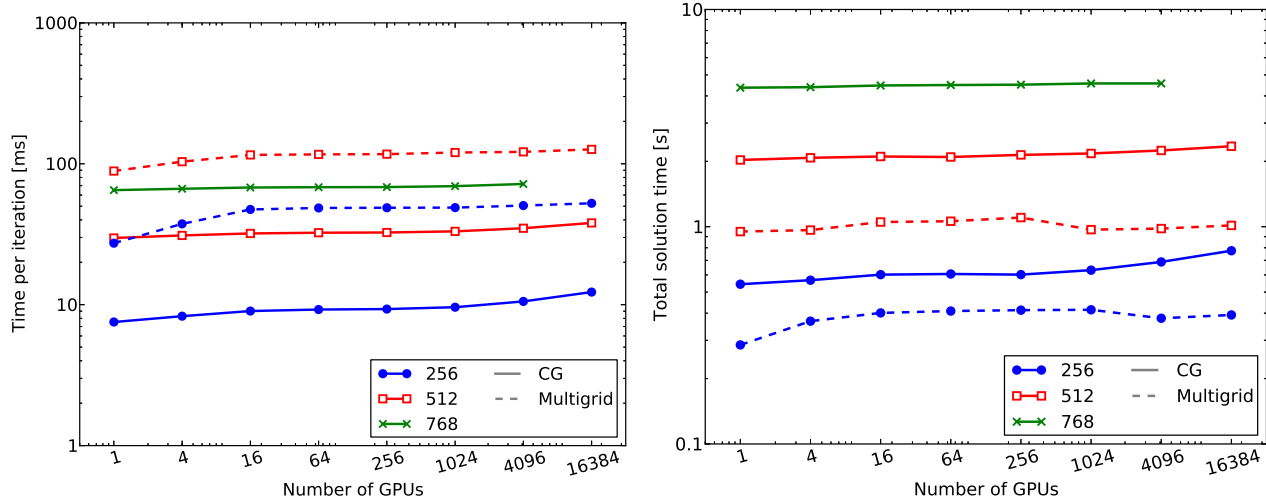


Figure 4: Weak scaling of the time per iteration (left) and total solution time (right) on different numbers of GPUs on Titan.

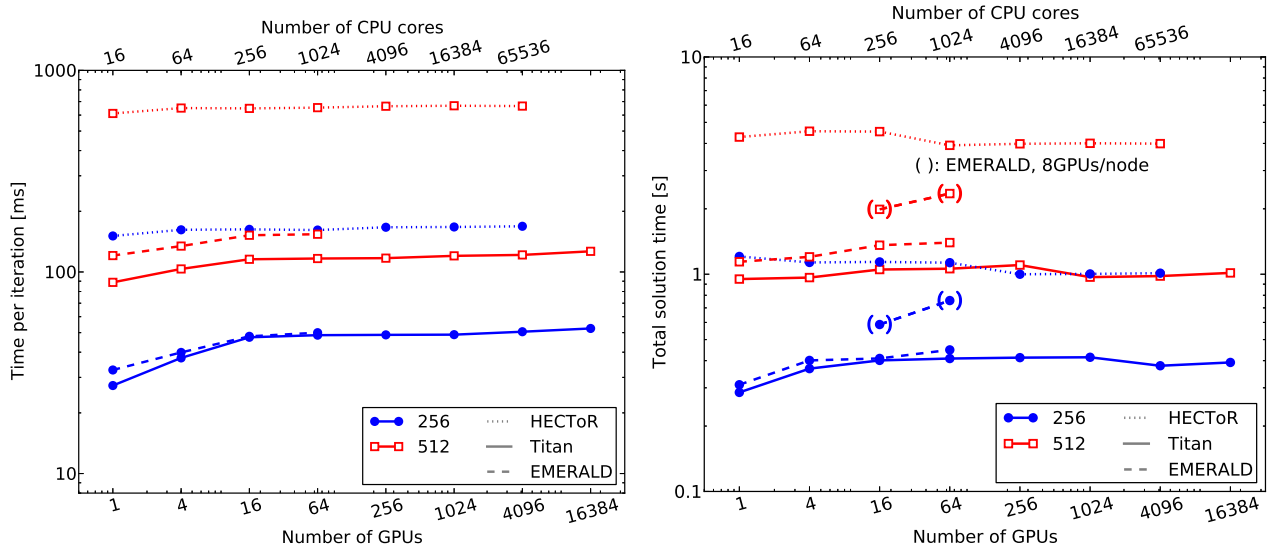


Figure 5: Time per iteration (left) and total solution time (right) of the multigrid solver on different clusters. Results obtained on EMERALD with 8 GPUs per node are shown in brackets.

# GPUs	# columns per GPU	CG		Multigrid	
		$t_{\text{solve}}[s]$	parallel efficiency	$t_{\text{solve}}[s]$	parallel efficiency
256	$512 \times 512 = 262,144$	2.141	—	1.102	—
1024	$256 \times 256 = 65,536$	0.631	84.8%	0.415	66.4%
4096	$128 \times 128 = 16,384$	0.541	24.8%	0.250	27.6%

Table 7: Strong scaling of the total runtime for a problem of size $8192 \times 8192 \times 128$ on Titan. The parallel efficiency is given relative to the 256 GPU run.

solver for two problem sizes on all three computer systems. On the EMERALD cluster two memory links have to be shared between all GPUs on a node and we found that we can achieve the best performance if we only use two GPUs per node, all results shown in Figure 5, except for the ones which are specifically marked, are obtained with this configuration. In summary the GPU implementation on Titan is roughly a factor 4 faster than the CPU implementation on HECToR (comparing one K20X GPU card to one 16 core AMD CPU) or, put differently, it is possible to solve a four times larger problem in the same total runtime. For the 512 problem we find that the CG solver on the GPU is about 6 times faster than the CPU implementation (not shown here).

The absolute performance on Titan both in terms of the global memory bandwidth and in terms of floating point operations per second was also quantified as described in Section 5.1. The absolute floating point performance is plotted in Figure 6 (left) for different problem sizes and numbers of GPUs. On one GPU the CG solver on the 512 problem can utilise about 5 % of the peak performance and on the largest core count it can still use 3 % of the peak FLOP rate of the entire machine. As both algorithms are memory bound, a more meaningful measure for the performance of the solvers is the global memory bandwidth. According to Figure 6 the CG solver can utilise between 30% and 60% of the theoretical peak global memory bandwidth. For the multigrid solver this number is smaller and in the range 15% - 25% on 16384 GPUs and 25% - 35% on one GPU. We stress again that we measure the “useful bandwidth” and obtaining a value close to 100% is very hard to achieve in practice.

It is also interesting to compare the performance on the two GPU systems. Although the difference in floating point performance between the cards is a factor two, one multigrid iteration on the Fermi card is only 20% – 35% slower than on the Kepler GPU, which is closer to what is expected due to the ratio between the global memory bandwidths on the two different cards ($BW_{\text{mem}}(K20X)/BW_{\text{mem}}(M2090) = (250 \text{ GByte/s})/(177 \text{ GByte/s}) \approx 1.4$). For the same number of total GPUs, using all 8 GPUs on an EMERALD node (instead of only using 2 per node) leads to a significant increase in the total runtime as all cards have to share the host-device bandwidth (see the additional data in Figure 5). This is expected as for the multigrid solver the cost of copying the right hand side to the GPU and copying the final solution back to the CPU is not negligible (see Table 3).

While we have not attempted to carry out detailed strong scaling runs, we report in Tab. 7 some preliminary results

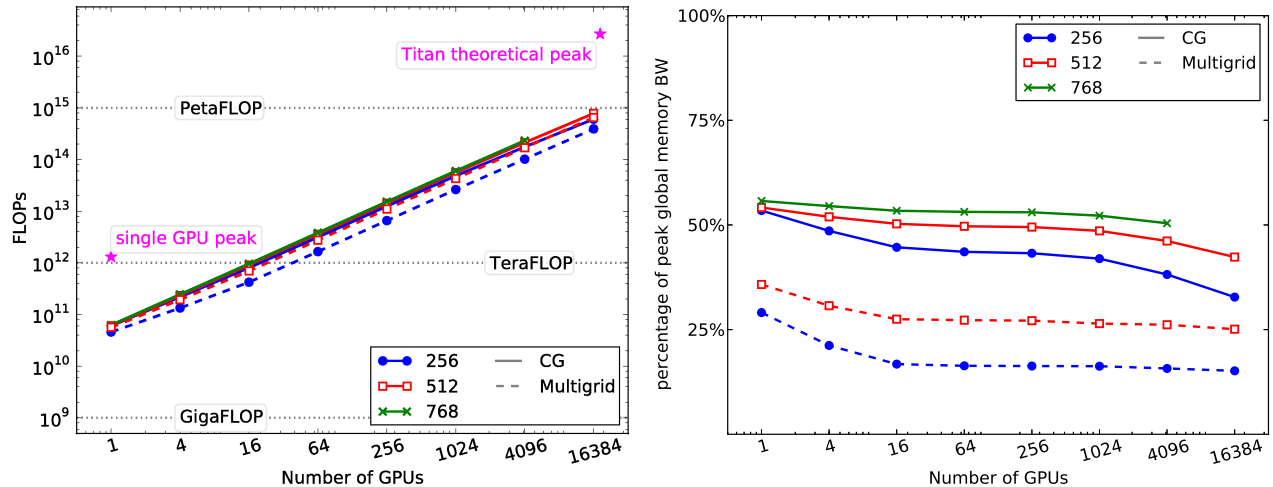


Figure 6: Floating point performance and percentage of utilised peak global memory bandwidth for different numbers of GPUs and local problem sizes n_x on Titan.

for a problem of size $8192 \times 8192 \times 128 = 8.6 \cdot 10^9$ unknowns on 256, 1024 and 4096 GPUs; the results were obtained on Titan. The efficiency is still around 85% for CG and 66% for the multigrid solver when the number of GPUs is increased from 256 to 1024, i.e. when the number of local columns on one GPU has gone from 512×512 to 256×256 . When the number of GPUs is increased further to 4096, the efficiency (measured relative to the 256 GPU run) drops to around 25% for both algorithms. On 4096 GPUs the number of threads per GPU is $128 \times 128 = 16384$. This has to be compared to the number of CUDA cores on the Kepler card which is 2688, i.e. each core will process only around 6 columns. As the local problem size is reduced further the number of columns per core will drop below 1 and it becomes impossible to exploit the full capacity of the GPU with our current implementation.

6. Conclusion and outlook

In this article we described massively parallel and efficient CUDA-C implementations of two memory bound iterative solvers for anisotropic PDEs. Equations of this type are encountered in many areas of geophysical modelling and we focus on the pressure correction equation which arises from semi-implicit semi-Lagrangian time stepping in atmospheric forecast models. The biggest gains can be achieved by choosing the algorithmically most efficient solver tailored to the problem, which in this case is a tensor-product geometric multigrid solver. It is about twice as fast as a preconditioned Conjugate Gradient method. We demonstrated the excellent absolute performance (measured in terms of the achieved memory bandwidth) of our solvers on Kepler GK110 GPUs and showed the very good weak scaling on up to 16384 GPUs of the Titan supercomputer. The GPU implementation is about a factor four faster than the CPU code on HECToR.

Although we believe that our implementation is close to optimal there are still potential improvements that could be considered. In this article we only give preliminary results for strong scaling of our solvers. In the strong scaling limit the number of threads will eventually be too small to utilise the full computational power of the GPU. This could be improved by exposing more parallelism in the algorithm, for example by using a parallel tridiagonal solver as described below. Furthermore the overhead from parallel communications between GPUs can be “hidden” by overlapping the halo exchange with computations. For this the domain is split up into an interior part and a boundary. After the calculations have been completed on the boundary, an asynchronous halo exchange is posted and calculations are continued for the interior part of the domain. While on a CPU this approach is straightforward, on a GPU the calculations on the boundary degrees of freedom will only be efficient if it is large enough in the x -direction to allow contiguous memory access in this direction. This is a problem in particular for the multigrid method on the coarse levels which are already very small. In some preliminary experiments we were not able to achieve any speedups with this technique.

Part of the reason for the poorer computational efficiency of the multigrid method is that the grids on the coarse levels are so small that the GPU can not be fully utilised, as we only parallelise in the horizontal direction due to the dependency in the tridiagonal solver. The amount of parallelism can be increased if we can assign several threads to work on each vertical column, by using a parallel tridiagonal solver such as cyclic reduction or a substructuring algorithm (see also [42]). Preliminary experiments with the substructuring method have shown that some gains with speedups of up to a factor two can be achieved on the coarser multigrid levels.

Currently our entire solver is implemented on the GPU, which means that the host CPU is idle. This is wasteful and could potentially be improved by splitting the work between the two processors. The obvious way of doing this is via splitting the domain between the CPU and GPU, as described for a shallow water solver in [33]. For the multigrid solver an alternative approach might be to process the finer levels on the GPU and the coarse levels on the CPU, thus minimising the amount of data that is copied between host and device. This could for example be done by using an additive multigrid algorithm, and potential benefits have to be balanced against worse algorithmic performance of the method.

Acknowledgements

We are grateful to Benson Muite for his help with porting and running the code on the Titan supercomputer and for kindly making part of his compute time allocation on the machine available for this project. We would like to thank Mauro Bianco (CSCS, Switzerland) for his help with using the GCL library and Mike Giles and István Reguly (Oxford) for useful discussions.

This work was funded as part of the NERC project on “Next Generation Weather and Climate Prediction” (NGWCP), grant number NE/K006762/1 and was supported also by European Regional Development Fund through the Estonian Centre of Excellence in Computer Science and the Estonian Science Foundation grant 9019. We used resources of the Oak Ridge Leadership Computing Facility at the Oak Ridge National Laboratory, which is supported by the Office of Science of the U.S. Department of Energy under Contract No. DE-AC05-00OR22725; the runs on Titan were carried out under DD Project CSC113. In addition we made use of the facilities of HECToR, the UK’s national high-performance computing service, which is provided by UoE HPCx Ltd at the University of Edinburgh, Cray Inc and NAG Ltd, and funded by the Office of Science and Technology through EPSRC’s High End Computing Programme. We would also like to acknowledge use of the EMERALD High Performance Computing facility provided via the Centre for Innovation (Cfi). The Cfi is formed from the universities of Bristol, Oxford, Southampton and UCL in partnership with STFC Rutherford Appleton Laboratory.

- [1] Eike Müller and Robert Scheichl. Massively parallel solvers for elliptic PDEs in numerical weather- and climate prediction. *accepted for publication in Q. J. Roy. Meteor. Soc.*, 2014.
- [2] Stefan Börm and Ralf Hiptmair. Analysis of tensor product multigrid. *Numer. Algorithms*, 26:200–1, 1999.
- [3] Markus Blatt and Peter Bastian. The Iterative Solver Template Library. In *Lecture Notes in Computer Science*, volume 4699, pages 666–675. Springer-Verlag, Berlin, Heidelberg, 2007.
- [4] Markus Blatt. A Parallel Algebraic Multigrid Method for Elliptic Problems with Highly Discontinuous Coefficients (PhD thesis). *Heidelberg*, 2010.
- [5] Van Emden Henson and Ulrike Meier Yang. BoomerAMG: a Parallel Algebraic Multigrid Solver and Preconditioner. *Applied Numerical Mathematics*, 41:155–177, 2000.
- [6] Hans Meuer, Erich Strohmaier, Jack Dongarra, and Horst Simon. Top500 list of supercomputers, November 2013. <http://www.top500.org/lists/2014/06/>. Accessed: 13 Aug 2014.
- [7] nVidia Corporation. nVidia Tesla GPU Accelerators data sheet. <http://www.nvidia.com/object/tesla-servers.html>. Accessed: 8 Jan 2014.
- [8] William C. Skamarock, Piotr K. Smolarkiewicz, and Joe B. Klemp. Preconditioned conjugate-residual solvers for Helmholtz equations in nonhydrostatic models. *Mon. Weather Rev.*, 125(4):587–599, 1997.
- [9] Stephen J. Thomas, Andrei V. Malevsky, Michel Desgagné, R. Benoit, P. Pellerin, and Michel Valin. Massively parallel implementation of the mesoscale compressible community model. *Parallel Comput.*, 23:2143 – 2160, 1997.
- [10] Abdessamad Qaddouri and Jean Côté. Preconditioning for an Iterative Elliptic Solver on a Vector Processor. In Jos Palma, A. Sousa, Jack Dongarra, and Vicente Hernandez, editors, *High Performance Computing for Computational Science VECPAR 2002*, volume 2565 of *Lecture Notes in Computer Science*, pages 451–455. Springer, Berlin, 2003.
- [11] Terry Davies, Mike J. P. Cullen, Andrew J. Malcolm, M. H. Mawson, Andrew Staniforth, A. A. White, and Nigel Wood. A new dynamical core for the Met Office’s global and regional modelling of the atmosphere. *Q. J. Roy. Meteor. Soc.*, 131(608):1759–1782, 2005.
- [12] William H. Press, Saul A. Teukolsky, William T. Vetterling, and Brian P. Flannery. *Numerical Recipes: The Art of Scientific Computing, 3rd Edition*. Cambridge University Press, New York, 2007.

- [13] Eike Müller, Xu Guo, Robert Scheichl, and Sinan Shi. Matrix-free GPU implementation of a preconditioned Conjugate Gradient solver for anisotropic elliptic PDEs. *to appear in Computation and Visualization in Science*, Feb 2014.
- [14] Mauro Bianco. An interface for halo exchange pattern. <http://www.prace-ri.eu/IMG/pdf/wp86.pdf>, 2013. Accessed: 11 Jan 2014.
- [15] Andreas Dedner, Eike Müller, and Robert Scheichl. Efficient multigrid preconditioners for atmospheric flow simulations at high aspect ratio. *International Journal for Numerical Methods in Fluids*, pages n/a–n/a, 2015. available online.
- [16] Jeff Bolz, Ian Farmer, Eitan Grinspun, and Peter Schröder. Sparse Matrix Solvers on the GPU: Conjugate Gradients and Multigrid. *ACM Transactions on Graphics*, 22:917–924, 2003.
- [17] Nolan Goodnight, Cliff Woolley, Gregory Lewin, David Luebke, and Greg Humphreys. A multigrid solver for boundary value problems using programmable graphics hardware. In *ACM SIGGRAPH 2005 Courses*, SIGGRAPH '05, New York, NY, USA, 2005. ACM.
- [18] Sandeep Menon and J. Blair Perot. Implementation of an efficient conjugate gradient algorithm for Poisson solutions on graphics processors. In *Proceedings of the 2007 Meeting of the Canadian CFD Society, Toronto Canada*, 2007.
- [19] Ali Cevahir, Akira Nukada, and Satoshi Matsuoka. Fast conjugate gradients with multiple GPUs. In *Computational Science–ICCS 2009*, pages 893–903. Springer Verlag, Berlin, Heidelberg, 2009.
- [20] Marco Ament, Günter Knittel, Daniel Weiskopf, and Wolfgang Strasser. A Parallel Preconditioned Conjugate Gradient Solver for the Poisson Problem on a Multi-GPU Platform. In *Parallel, Distributed and Network-Based Processing (PDP), 2010, 18th Euromicro International Conference on*, pages 583–592, feb. 2010.
- [21] Ali Cevahir, Akira Nukada, and Satoshi Matsuoka. High performance conjugate gradient solver on multi-GPU clusters using hypergraph partitioning. *Computer Science-Research and Development*, 25(1-2):83–91, 2010.
- [22] Serban Georgescu and Hiroshi Okuda. Conjugate gradients on multiple GPUs. *International Journal for Numerical Methods in Fluids*, 64(10-12):1254–1273, 2010.
- [23] Michael Griebel and Peter Zaspel. A multi-GPU accelerated solver for the three-dimensional two-phase incompressible Navier-Stokes equations. *Computer Science-Research and Development*, 25(1-2):65–73, 2010.
- [24] Dana A Jacobsen, Julien C Thibault, and Inanc Senocak. An MPI-CUDA implementation for massively parallel incompressible flow computations on multi-GPU clusters. In *48th AIAA Aerospace Sciences Meeting and Exhibit*, volume 16, 2010.
- [25] Zhuo Feng and Zhiyu Zeng. Parallel multigrid preconditioning on graphics processing units (GPUs) for robust power grid analysis. In *Proceedings of the 47th Design Automation Conference*, pages 661–666. ACM, 2010.
- [26] Markus Geveler, Dirk Ribbrock, Dominik Göttsche, Peter Zajac, and Stefan Turek. *Efficient finite element geometric multigrid solvers for unstructured grids on GPUs*. Techn. Univ., Fak. für Mathematik, 2011.
- [27] Dana A. Jacobsen and Inanc Senocak. A full-depth amalgamated parallel 3D geometric multigrid solver for GPU clusters. In *49th AIAA Aerospace Science Meeting*, 2011.
- [28] Nathan Bell and Michael Garland. Efficient sparse matrix-vector multiplication on CUDA. NVIDIA Technical Report NVR-2008-004, NVIDIA Corporation, December 2008.
- [29] Timothy A. Davis and Yifan Hu. The University of Florida Sparse Matrix Collection. *ACM Trans. Math. Softw.*, 38(1):1:1–1:25, December 2011.
- [30] Gundolf Haase, Manfred Liebmann, Craig C Douglas, and Gernot Plank. A parallel algebraic multigrid solver on graphics processing units. In *High performance computing and applications*, pages 38–47. Springer, 2010.
- [31] James Brannick, Yao Chen, Xiaozhe Hu, and Ludmil Zikatanov. Parallel Unsmoothed Aggregation Algebraic Multigrid Algorithms on GPUs. In Oleg P. Iliev, Svetozar D. Margenov, Peter D Minev, Panayot S. Vassilevski, and Ludmil T Zikatanov, editors, *Numerical Solution of Partial Differential Equations: Theory, Algorithms, and Their Applications*, volume 45 of *Springer Proceedings in Mathematics & Statistics*, pages 81–102. Springer New York, 2013.

- [32] Dominik Göttsche, Robert Strzodka, Jamaludin Mohd-Yusof, Patrick McCormick, Hilmar Wobker, Christian Becker, and Stefan Turek. Using GPUs to improve multigrid solver performance on a cluster. *International Journal of Computational Science and Engineering (IJCSE)*, 4(1):36–55, 2008.
- [33] Chao Yang, Wei Xue, Haohuan Fu, Lin Gan, Linfeng Li, Yangtong Xu, Yutong Lu, Jiachang Sun, Guangwen Yang, and Weimin Zheng. A Peta-scalable CPU-GPU Algorithm for Global Atmospheric Simulations. *SIGPLAN Not.*, 48(8):1–12, February 2013.
- [34] Björn Gmeiner, Harald Köstler, Markus Stürmer, and Ulrich Rüde. Parallel multigrid on hierarchical hybrid grids: a performance study on current high performance computing clusters. *Concurrency and Computation: Practice and Experience*, 26(1):217–240, 2014.
- [35] Ulrich Trottenberg, Cornelis W. Oosterlee, and Anton Schüller. *Multigrid*. Academic Press, San Diego, London, Sydney, Tokyo, 2001.
- [36] Nigel Wood, Andrew Staniforth, Andrew White, Tom Allen, Michail Diamantakis, Markus Gross, Thomas Melvin, Chris Smith, Simon Vosper, Mohamed Zerroukat, and John Thuburn. An inherently mass-conserving semi-implicit semi-Lagrangian discretisation of the deep-atmosphere global nonhydrostatic equations. *Q. J. Roy. Meteor. Soc.*, 2013. Published online December 4th 2013.
- [37] Robert Sadourny. Conservative Finite-Difference Approximations of the Primitive Equations on Quasi-Uniform Spherical Grids. *Mon. Weather Rev.*, 100(2):136–144, 1972.
- [38] Yousef Saad. *Iterative Methods for Sparse Linear Systems, Second Edition*. Society for Industrial and Applied Mathematics, Philadelphia, 2003.
- [39] William. L. Briggs, Van Emden Henson, and Steve F. McCormick. *A Multigrid Tutorial*. Society for Industrial and Applied Mathematics, Philadelphia, 2000.
- [40] AMD Corporation. . <http://www.amd.com/uk/products/server/processors/6000-series-platform/6200/Pages/6200-series-processors.aspx>. Accessed: 12 Jan 2014.
- [41] Mark Harris. An Efficient Matrix Transpose in CUDA C/C++. <http://devblogs.nvidia.com/parallelforall/efficient-matrix-transpose-cuda-cc/>. Accessed: 7 Jan 2014.
- [42] Endre Laszlo. Efficient Solution of Multiple Scalar and Block-Tridiagonal Equations. GPU Technology Conference, 2014.