

An Accelerated 3D Navier–Stokes Solver for Flows in Turbomachines

Tobias Brandvik
e-mail: tb302@cam.ac.uk

Graham Pullan¹
e-mail: gp10006@cam.ac.uk

Department of Engineering,
Whittle Laboratory,
University of Cambridge,
1 JJ Thomson Avenue,
Cambridge CB3 0DY, UK

A new three-dimensional Navier–Stokes solver for flows in turbomachines has been developed. The new solver is based on the latest version of the Denton codes but has been implemented to run on graphics processing units (GPUs) instead of the traditional central processing unit. The change in processor enables an order-of-magnitude reduction in run-time due to the higher performance of the GPU. The scaling results for a 16 node GPU cluster are also presented, showing almost linear scaling for typical turbomachinery cases. For validation purposes, a test case consisting of a three-stage turbine with complete hub and casing leakage paths is described. Good agreement is obtained with previously published experimental results. The simulation runs in less than 10 min on a cluster with four GPUs. [DOI: 10.1115/1.4001192]

1 Introduction

A key metric in the evaluation of a computational fluid dynamics (CFD) solver is the time taken per node per timestep. Advances in hardware and, to a lesser extent, algorithms have enabled this metric to fall continuously. In terms of low-cost commodity hardware [central processing units (CPUs)], processor clock speed and changes in processor architecture were the primary drivers for this advance in the 1980s and 1990s. At present, parallel computing using clusters of multicore processors is the key enabler. Nonetheless, a step-change in this metric (at least one order-of-magnitude) would be invaluable in bringing high-fidelity [large eddy simulation (LES) and direct numerical simulation (DNS)] solutions into routine industrial use or in making the current standard of design tools interactive.

Alongside the shift from single- to multicore CPUs that has occurred over the last 5 years, advances have also been made for other types of processors. There are now a variety of chips on the market that exhibit a higher level of parallelism (and hence performance) than CPUs, including graphics processing units (GPUs) and the Sony/Toshiba/IBM (STI) cell microprocessor. Taking advantage of such processors for CFD calculations could deliver a step-change in performance today but requires significant changes to the underlying code.

There has been only limited work reported so far on the use of these novel processors for CFD. The present authors have presented results for 2D [1] and 3D [2] Euler solvers for turbomachinery applications, achieving speed-ups of an order-of-magnitude for both AMD and NVIDIA GPUs compared with a single CPU core. Similar results, with the extension to include a full multigrid scheme, have also been presented by Elsen et al. [3] for a 3D Euler solver with applications to external flows.

In this paper, a new three-dimensional Navier–Stokes solver, which runs on GPUs is presented. The solver, called TURBOSTREAM, includes a mixing-length turbulence model and the capability of simulating both steady and unsteady flows in multirow turbomachines. In addition, the solver is able, through the use of the message passing interface (MPI), to utilize many GPUs together to solve problems that are beyond the scope of a single-processor. To the authors' knowledge, this represents the first vis-

cous solver for engineering flows running on a large number of GPUs.

We first present an overview of many-core processors, focusing specifically on GPUs and how they relate to CPUs. A brief description of the solver algorithm and implementation is then given, followed by a discussion of the solver's performance as it compares to an older CPU solver that implements the same algorithm. Finally, a test case consisting of a three-stage turbine with hub and casing leakage paths is presented. Comparisons showing good agreement with experiments are also given.

2 Many-Core Processors

In order to keep increasing the computational power of their chips, the semiconductor industry has now moved from single- to multicore designs. The reason for this switch is the requirement to keep within an acceptable power envelope of around 100–200 W. Given this constraint, Borkar [4] identified the diminishing returns of core complexity as the main motivator for multicore processors. He referred to this relationship as Pollack's rule, which states that the computational power of a core is roughly proportional to the square root of its complexity. Therefore, it is clearly more power-efficient to use the extra transistors offered by technology scaling to add more cores to a chip, rather than to increase the complexity of the existing ones. This trend is widely expected to continue, resulting in 1000 core chips being common-place within the next 10 years.

For software development, the consequences of this explosion in core count are far-reaching, requiring significant changes to old codes and the algorithms they use. In this paper, the problem is approached from the point of view of a developer rewriting an existing Fortran structured grid CFD solver. To set the stage, we first introduce the two different processors that will be considered. The first is a quad-core Intel Xeon CPU, which is representative of the processors that run most CFD solvers today; the second is the less familiar design of an NVIDIA GPU that is the current target processor of TURBOSTREAM. A schematic overview of both processors can be seen in Fig. 1.

2.1 Intel Xeon. The processor considered here is the 2.33GHz Harpertown variant in Intel's Xeon line. It is a general-purpose processor, meaning that it is capable of running an operating system on its own. Harpertown is a dual-die quad-core, i.e., two dual-cores put together in the same package. Each core has its own 32 KB L1 cache while each die has its own 4 MB L2 cache shared between the two cores. Each core also has an adder and a multiplier for 128 bit vectors, making it capable of 18.6 single

¹Corresponding author.

Contributed by the International Gas Turbine Institute (IGTI) of ASME for publication in the JOURNAL OF TURBOMACHINERY. Manuscript received August 3, 2009; final manuscript received December 21, 2009; published online October 27, 2010. Editor: David Wisler.

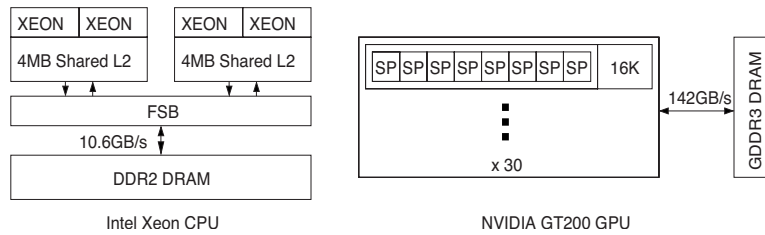


Fig. 1 CPU and GPU architecture overview

precision GFLOP/s (10^9 floating point operations per second). The theoretical aggregate performance of the whole chip is therefore 74.4 GFLOP/s. The cores have access to external memory through the front side bus at a rate of 10.6 GB/s.

The programming approach used for structured grid applications on CPUs is now well established. Typically, some form of domain decomposition with ghost cells is used to split the domain between the processor cores. One CPU process is then started per core. Each process iterates over its part of the domain, updating the grid variables as it goes along. At the end of each iteration, the processes exchange data with each other to update the ghost cells by using MPI. A variation in this approach is to only start one CPU process per processor, and then parallelize the work given to that processor across its cores using threads (either created explicitly or through compiler directives such as OpenMP). However, since most established codes were already parallelized with MPI before the arrival of multicore CPUs, the pure MPI approach involves less work and seems to be more popular.

2.2 NVIDIA GPU. The NVIDIA GPU considered here is the latest GT200 chip. The older G80 chip is also used in some of the performance measurements presented later—this has approximately half the performance of the GT200.

The GT200 is designed to accelerate the rendering of 3D scenes in computer games, so its volume sales are driven by the computer games industry. Unlike a CPU, it is not a general-purpose chip and cannot run an operating system. Instead, it is sold as part of an add-in card (graphics card) that comes with its own on-board memory and plugs into an expansion slot on the PCI-express bus. The GT200 consists of 30 multiprocessors (MP), each of which contains eight scalar processing units (SP) and 16 KB of explicitly managed local storage (referred to as shared memory). Each MP has its own instruction counter and operates independently of the others. Each SP can schedule one multiply and one multiply-add operation (both single precision) per cycle, giving a theoretical peak performance of 933GFLOP/s at 1.296GHz. By using a wide 512 bit bus to the graphics card's on-board GDDR3 memory, the GT200 achieves a maximum bandwidth of 141.7GB/s.

To simplify the programming of their GPUs, NVIDIA has developed an extension to the C programming language called CUDA. In a CUDA program, the developer sets up a large number of threads (often several thousand) that are grouped into thread blocks. A CUDA thread is the smallest unit of execution and has a set of registers and a program counter associated with it. This is similar to traditional CPU threads but CUDA threads are much less expensive to create and swap between. Each thread block is executed on a single multiprocessor. It is possible to synchronize the threads within a block, allowing the threads to share data through the shared memory. Given that a thread block can consist of more threads than the number of processors in a multiprocessor, the hardware is responsible for scheduling the threads. This allows it to hide the latency of fetches from the on-board memory by letting some threads perform computations while others wait for data to arrive. For structured grid applications, a natural way of organizing the code is to split the main grid into smaller grids,

which can fit into the shared memory. A block of threads is then started within each of the smaller grids to compute the updated variables.

3 Algorithm

TURBOSTREAM is heavily based on the long line of codes from Denton. In particular, it uses the same algorithm as that of the latest Denton code, TBLOCK, with only minor differences in the way that this algorithm is implemented. A complete description of TBLOCK is given by Klostermeier [5] while shorter overviews and examples of its application to turbomachinery problems were published by Reid et al. [6] and Rosic et al. [7]. In addition, the motivation for the current method can be traced through a series of papers by Denton [8–11]. Here, we give a basic overview of the algorithm as it is used in TURBOSTREAM.

3.1 Algorithm Overview. TURBOSTREAM uses a multiblock topology with arbitrary patch interfaces to capture complex geometries. Information is passed between blocks using surface patches, which contain nodes that are physically coincident but reside on different blocks. The flow properties at these nodes are averaged at the end of each timestep. Parallel simulations on a cluster of processors can be performed by decomposing the domain on a block basis. This decomposition is performed as an automatic preprocessing step in which each block can be further split into smaller blocks to achieve better load-balancing.

The Navier–Stokes equations are discretized using a finite volume method with vertex storage in a structured grid of hexahedral cells. In this technique, the equations in their integral form for mass, momentum, and energy are used.

Mass

$$\frac{\partial}{\partial t} \int_{\Omega} \rho d\Omega + \oint_A \rho \mathbf{u} \cdot d\mathbf{A} = 0 \quad (1)$$

Momentum

$$\frac{\partial}{\partial t} \int_{\Omega} \rho \mathbf{u} d\Omega + \oint_A \rho \mathbf{u} (\mathbf{u} \cdot d\mathbf{A}) + \oint_A p d\mathbf{A} - \oint_A \boldsymbol{\tau} \cdot d\mathbf{A} = 0 \quad (2)$$

Energy

$$\frac{\partial}{\partial t} \int_{\Omega} \rho e d\Omega + \oint_A \rho h_0 \mathbf{u} \cdot d\mathbf{A} - \oint_A (\boldsymbol{\tau} \cdot \mathbf{u}) \cdot d\mathbf{A} - \oint_A \lambda \nabla T \cdot d\mathbf{A} = 0 \quad (3)$$

where Ω is a control volume bounded by a surface A . The above integrals are performed on each cell in the grid using a second-order spatial discretization.

The equations can be expressed less formally as

$$\frac{\partial \mathbf{U}}{\partial t} = \frac{\sum \mathbf{F}}{V} + \mathbf{S} \quad (4)$$

where \mathbf{U} is a vector containing the primary flow variables, \mathbf{F} is a vector containing the fluxes of the primary variables, the flux

summation is over the faces of the cell, \mathbf{S} contains any source terms, and V is the volume of the cell. In the manner described by Denton [10], the source vector is here used to hold the viscous terms. This equation is integrated forward to reach a steady-state ($d\mathbf{U}/dt \approx 0$) using the Scree scheme (see Denton [12]).

$$\Delta \mathbf{U} = \left(2 \frac{\partial \mathbf{U}}{\partial t} \Big|_n - \frac{\partial \mathbf{U}}{\partial t} \Big|_{n-1} \right) \Delta t \quad (5)$$

where the subscripts refer to the timestep that the derivatives were evaluated at.

Since the integrals are evaluated for a hexahedral cell and the solver uses vertex storage, the cell-based $\Delta \mathbf{U}$ has to be distributed to the surrounding vertices, each receiving an eighth. Finally, to maintain numerical stability, artificial smoothing is then applied to all the flow variables. Presently, only second-order smoothing is used but a more traditional blended second- and fourth-order smoothing procedure is being considered.

3.2 Turbulence Model. The effect of turbulence is modeled using a simple algebraic mixing-length model in which the turbulent viscosity ν_t is related to length scale l_{mix} over which turbulent mixing is assumed to take place.

$$\nu_t = l_{\text{mix}}^2 \sqrt{2S_{ij}S_{ij}} \quad (6)$$

where S_{ij} is the strain-rate tensor

$$S_{ij} = \frac{1}{2} \left(\frac{\partial u_i}{\partial x_j} + \frac{\partial u_j}{\partial x_i} \right) \quad (7)$$

The main drawback of this model is the specification of the mixing-length, which is different for every type of flow. Experience from previous Denton codes has shown that for turbomachinery applications, a limiter based on the blade pitch is appropriate.

$$l_{\text{mix}} = \begin{cases} \kappa y_n, & y_n < x_{\text{lim}} \\ \kappa x_{\text{lim}}, & y_n > x_{\text{lim}} \end{cases} \quad (8)$$

where κ is a constant, y_n is the normal distance from the nearest wall, and x_{lim} is usually taken to be 3% of the pitch. y_n is calculated by TURBOSTREAM before the start of the main timestepping loop using the Poisson equation approach described by Tucker et al. [13].

To avoid having to use many grid points in the boundary layer, the flow is allowed to slip at the walls and a wall-function is used to obtain an expression for the wall shear stress. In this approach, it is assumed that the first grid point away from the wall lies either in the viscous sublayer or in the logarithmic region of a turbulent boundary layer. In the former case the wall shear stress is approximated by

$$C_{f,w} = \frac{1}{\text{Re}_w} \quad (9)$$

and in the latter case by a curve fit to the log-law in the form of

$$C_{f,w} = -0.001767 + \frac{0.03177}{\ln \text{Re}_w} + \frac{0.25614}{(\ln \text{Re}_w)^2} \quad (10)$$

where $C_{f,w}$ is the coefficient of friction defined as

$$C_{f,w} = \frac{\tau_w}{\frac{1}{2} \rho U_w^2} \quad (11)$$

and Re_w is the cell Reynolds number defined as

$$\text{Re}_w = \frac{\rho U_w y_w}{\mu} \quad (12)$$

In the above equations, U_w is the velocity at the first grid node off the wall and y_w is the height of the cell normal to the wall.

3.3 Convergence Acceleration. To accelerate the convergence rate, TURBOSTREAM uses both spatially varying timesteps and a multigrid scheme. In the former method, the timestep in each cell is limited by the local flow properties and geometry, allowing much larger timesteps to be used in the large free-stream cells that would otherwise be limited by the small cells in the boundary layer. The latter method uses multiple grid levels, each coarser than the preceding one, to accelerate the convergence by dispersing transients quickly on the coarser levels while retaining the spatial accuracy of the finest. In the Denton formulation used by TURBOSTREAM, adjacent cells are combined to form a grid of larger cells or blocks. The new coarse mesh is treated in just the same way as the original fine grid and so much larger timesteps are possible. The change in the value of \mathbf{U} during one iteration is given by

$$\Delta \mathbf{U} = \left[\frac{\Sigma \mathbf{F}}{V} + \mathbf{S} \right]_{\text{cell}} \Delta t_{\text{cell}} + \sum_{\text{blks}} \left[\left[\frac{\Sigma \mathbf{F}}{V} + \mathbf{S} \right]_{\text{block}} \Delta t_{\text{block}} \right] \quad (13)$$

where the summation is over all the blocks to which the relevant cell belongs. Typically, three levels of multigrid with a coarsening ratio of two are used.

3.4 Multistage and Unsteady Simulations. For steady-state calculations, multistage simulations are made possible by circumferentially averaging the flow entering a blade row, see Denton and co-worker [14,15,11]. The technique is applied at a mixing plane, the position of which is arbitrarily chosen by the user. The nonuniform flow upstream of this plane is mixed out so that it becomes pitchwise but not spanwise, uniform downstream of the mixing plane. The process is conservative but like any mixing process, it is irreversible (see Fritsch and Giles [16]). In addition, care must be taken to avoid locating the mixing plane too close to leading or trailing edges and, thereby, enforcing a nonphysical circumferentially uniform flow.

For unsteady simulations, TURBOSTREAM implements the “dual timestepping” technique proposed by Jameson [17]. This procedure allows the convergence acceleration methods described earlier to be used in time-accurate simulations by splitting the calculation up into a number of implicit “outer loops” that iterate forward in real time, each of which is comprised of a number of explicit “inner loops” that converge the flow to a steady-state in pseudotime. Multistage unsteady simulations are enabled through an interpolation procedure that transfers information across a sliding interface that connects the upstream and downstream grid blocks, which rotate relative to each other.

4 Implementation

During the initial considerations of implementing a flow solver to run on many-core architectures, there was some concern over the great number of different processors and programming models available. For GPUs alone, there are several options; these include the two main chip manufacturers AMD and NVIDIA and at least half a dozen programming languages and libraries that in some cases only target one vendor’s GPUs. In this situation, it seems a daunting prospect to pick one combination that will have the longevity required for CFD solvers whose lifetime is often measured in decades.

For this reason, it was decided that another approach than a direct implementation was needed. TURBOSTREAM is therefore expressed as a series of subroutine definitions in the high-level python scripting language. These definitions contain the input and output arguments for each subroutine, as well as the computations that are carried out within them. An as yet unpublished source-to-source compiler, which was developed for this work by the authors, is then used to transform these definitions into source code that can be further compiled for the target architecture that we wish to run on. Currently, the compiler can produce code for either multicore CPUs or NVIDIA GPUs, with support for the Cell processor currently being developed. Aside from enabling the

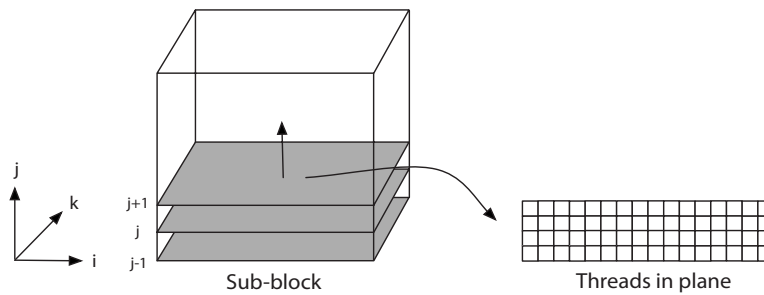


Fig. 2 Iteration procedure for stencil subroutines

solver to run on many different processors from the same source code definition, this source-to-source compilation approach has two other main benefits:

1. Since the definitions of the solver subroutines are completely separate from the source code that is actually produced for the target platform, the compiler is free to perform many different optimizations that would otherwise have complicated the code to an unacceptable degree. For an indication of the range and complexity of the many optimizations of which only a subset are currently used by our compiler, necessary to achieve near-optimal performance on modern many-core architectures, see Datta et al. [18].
2. The use of a high-level language to express the logic of the solver makes it easier for domain scientists to add extensions such as new numerical schemes and turbulence models. This capability is important because the solver is intended to be a platform for academic research as well as a production code for day-to-day turbomachinery design.

In addition to the compiler, a runtime library has also been developed that takes care of memory management, subroutine invocation, file input/output (I/O) and MPI communication. Since the details of the former two tasks are different for each processor, these parts of the library have to be written separately for each processor.

It should be noted that the source-to-source compilation strategy described here is only possible because of the limited range of computations that are performed by structured grid solvers. In short, each subroutine is a combination of stencil operations that use the nearest neighbors of a node to update its properties. The only computations that break with this paradigm are in the multi-grid routine, which therefore has to be implemented separately for each processor.

The main difficulty in producing efficient code for the stencil subroutines is in parallelizing the computations across the hundreds of scalar processing units present on modern GPUs. The strategy used by the compiler is to split each grid block into smaller sub-blocks that are computed independently from each other. One CUDA thread is started for each node in a plane of the sub-block. At the start of the subroutine, each thread loads in the necessary grid node values corresponding to its location. The threads then iterate upwards in the sub-block (Fig. 2), each time fetching a new plane and computing the values for the current one. Given that the shared memory can hold 16 KB, a typical sub-block size is $16 \times 10 \times 5$. For such sub-block sizes, the surface-to-volume ratio is low enough to get good data reuse. The overall approach is similar to that described by Williams et al. [19] for structured grid applications on the cell processor. A more detailed explanation of this implementation strategy and how it differs from that commonly used on CPUs has been included in Appendix.

A final issue that has to be considered for a GPU implementation is that of data transfers across the PCI-Express bus, which bridges the CPU and GPU memory spaces. The PCI-Express bus

has a theoretical maximum bandwidth of 4 GB/s or 8 GB/s depending on whether it is of generation 1 or 2. When this number is compared with the bandwidth between the GPU's on-board GDDR3 memory and the GPU multiprocessors (up to 141.7 GB/s), it becomes clear that any algorithm that requires a large amount of continuous data transfer between the CPU and GPU will not achieve good performance. For a CFD solver, the obvious solution is to limit the size of the domain that can be calculated so that all of the necessary data can be stored in the GPU's on-board memory. Using this approach, it is only necessary to perform large transfers across the PCI-Express bus at the start of the calculation (the geometry) and at the end (the final flow solution). High-end GPUs today have up to 4 GB of on-board memory, sufficient to store all the data needed by TURBOSTREAM for a grid with 12×10^6 nodes, so this restriction is not a significant limitation.

When operating in parallel across multiple GPUs, some boundary information must inevitably be transferred across the PCI-Express bus at the end of every timestep. However, as will be shown in the next section, the low surface-to-volume ratios in turbomachinery grids mean that this data transfer is not a bottleneck.

5 Performance

For any CFD solver, there are two important performance metrics.

1. How fast is the solver on a single-processor?
2. How well does the performance of the solver scale when multiple processors are used together to tackle larger problems?

TURBOSTREAM dramatically increases the single-processor performance as compared with other solvers by using NVIDIA GPUs instead of traditional CPUs. To demonstrate this speed-up, we compare TBLOCK running on all four cores of an Intel Xeon 2.33 GHz CPU with TURBOSTREAM running on an NVIDIA GT200 GPU. TBLOCK was compiled with the Intel 10.1 Fortran compilers with automatic vectorization and the highest degree of optimization turned on while TURBOSTREAM was compiled with NVIDIA's GPU compiler. Both solvers show approximately constant performance for grids with more than 10^5 nodes so a representative case with 10^6 nodes was used. The results are summarized in Table 1, which shows the time taken per grid node per timestep for each solver. The high performance of TURBOSTREAM running on the NVIDIA GPU is primarily due to the GPU's higher memory bandwidth, as well as the extra optimizations allowed by the source-

Table 1 Single processor performance

Solver	Processor	Time/node/step
TBLOCK	Intel Xeon 2.33 GHz	5.1×10^{-7} s
TURBOSTREAM	NVIDIA GT200	2.7×10^{-8} s

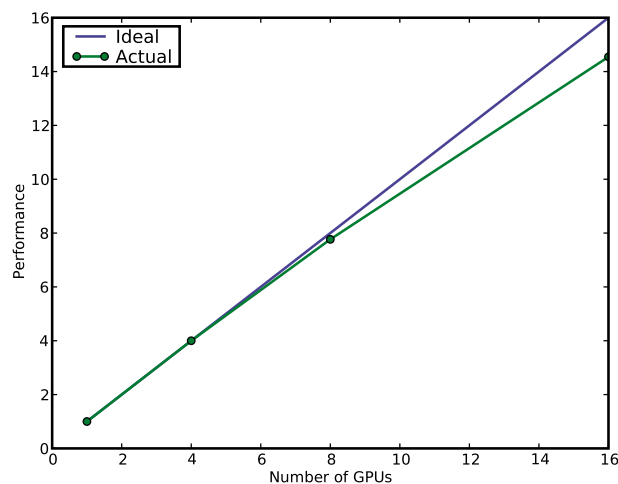


Fig. 3 TURBOSTREAM weak scaling over multiple GPUs. Performance is measured as the inverse of the time per grid node per timestep.

to-source compilation. In terms of wall-clock time, both solvers converge in approximately the same number of timesteps, with a typical 300,000 node single-row calculation that requires 5000 steps taking approximately 1 min in total with TURBOSTREAM and 20 min in total with TBLOCK.

Over the last 10 years, CFD has become increasingly reliant on clusters of processors to enable more detailed simulations within design time frames. For this reason, the scalability of a solver across multiple processors can be equally important as its single-processor performance. A potential problem with increasing the single-processor performance by an order-of-magnitude is then that the multiprocessor performance suffers since the time required to exchange boundary information remains roughly constant. However, the low surface-to-volume ratios in turbomachinery grids mean that good scalability can be achieved even with very fast solvers. To demonstrate this point, Fig. 3 shows the performance of TURBOSTREAM across a cluster of 16 NVIDIA G80 GPUs. There are four nodes in the cluster, each consisting of a traditional 1U server with a quad-core CPU connected through PCI-Express cables to another 1U server with four GPUs. The nodes are networked together with 1 Gigabit Ethernet interconnects.

To simplify the mesh generation for an arbitrary number of processors, we use an idealized case of simple flow through a square channel. In an attempt to represent a typical multistage turbomachinery calculation with one stage per GPU, the ratio between the number of points in the axial, radial and circumferential directions is taken to be 4:1:1. Two million nodes are used per GPU, so the total size of the simulation scales with the number of

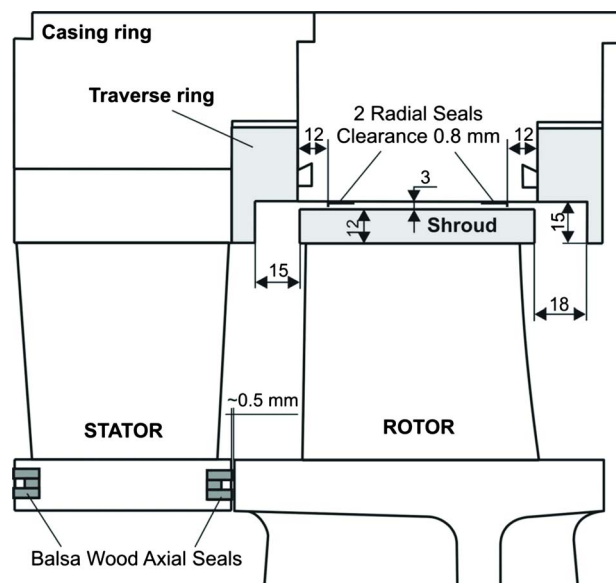


Fig. 4 Single stage geometry

GPUs used. In the authors' experience, this setup closely resembles real world usage—due to the high single-processor performance of the solver, multiple GPUs are only used in practice if the simulation is too large to fit in a single GPU's on-board memory.

As can be seen in Fig. 3, almost ideal scaling is obtained for 16 GPUs. It should also be noted that the interconnect used here (Gigabit Ethernet) has poor performance compared with other options, and that using a higher performance interconnect such as Infiniband should improve the scaling performance further.

6 Validation

Validation is the biggest hurdle for any new flow solver to gain acceptance in the community. The authors are currently running through many different existing TBLOCK test cases with TURBOSTREAM. Although minor differences between the implementation of the two solvers mean that the results are not always identical, they are in all cases in close agreement with each other. A calculation of a three-stage turbine with leakage paths is presented in this work.

6.1 Three-Stage Turbine With Leakage Paths. The test case is a three-stage turbine with leakage paths. It was originally presented by Rosic et al. [7] to demonstrate the importance of shroud leakage modeling in multistage turbine flow calculations. The original work was carried out using TBLOCK. Here, we show that TURBOSTREAM is capable of producing similar results.

The arrangement of a single stage is presented in Fig. 4, show-

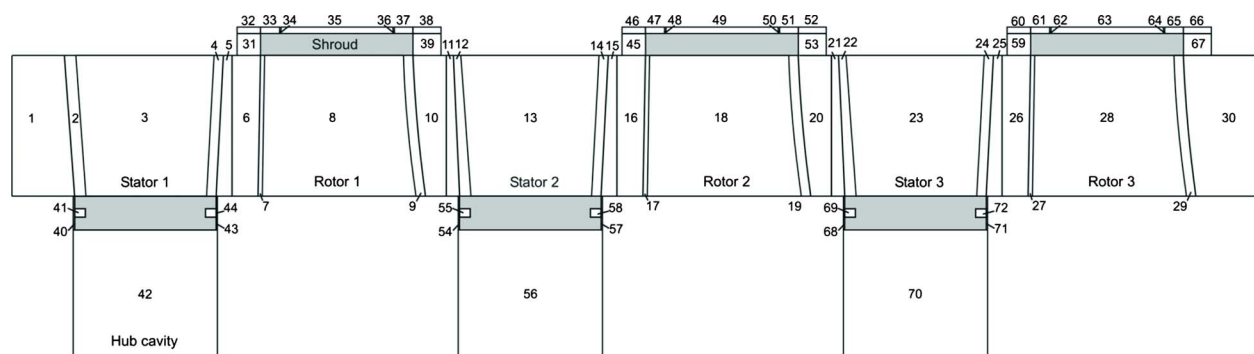


Fig. 5 Computational domain

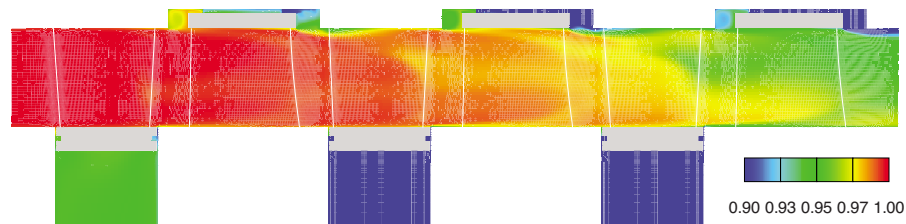


Fig. 6 Pitchwise averaged entropy function: $\exp(-\Delta s/R)$ (TURBOSTREAM)

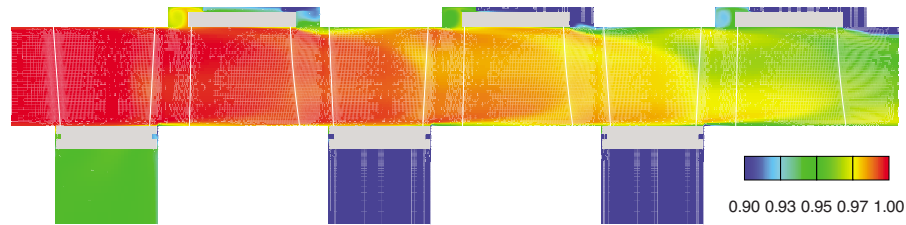


Fig. 7 Pitchwise averaged entropy function: $\exp(-\Delta s/R)$ (TBLOCK)

ing both the hub and shroud geometries. All leakage paths are fully represented in the CFD mesh and a cavity with rotating walls has been added to represent the area below the hub. This stage is replicated three times to form the whole machine, resulting in an overall computational domain as shown in Fig. 5. The mesh used is of the H-type and the total number of grid nodes is 4.5×10^6 . A cluster with four NVIDIA G80 GPUs was used to calculate the flow, resulting in an overall run-time of less than 10 min. Further computational and experimental details are given in the original paper.

Experimental and numerical results are presented using spanwise distributions of the pitchwise averaged exit yaw angle for the third stator and rotor, as well as exit total pressure coefficient contours for the third stator. The total pressure coefficient was obtained by nondimensionalizing the total pressure by the total pressure drop across the whole machine:

$$C_{p_0} = \frac{p_{0_{in}} - p_0}{p_{0_{in}} - p_{0_{ex}}} \quad (14)$$

Two sets of TURBOSTREAM results are presented: one with leakage flows and one without. For comparison, TBLOCK results for the case with leakages are also shown.

As should be expected, the results are similar to those of the original work (Rosic et al. [7]) and so only a brief discussion is warranted here. Figure 6 shows a meridional view of pitchwise averaged entropy function for TURBOSTREAM (Fig. 7 shows the same for TBLOCK). It is clear that the rotor shroud leakage flows interact strongly, and enhance the casing secondary flow in the following stator. Similarly, but to a lesser extent, the stator hub leakages add to the strength of the hub secondary flow in the following rotor. As an example of this effect, Fig. 8 shows a comparison of measured (using a five-hole pneumatic probe) and calculated total pressure loss coefficient at the exit of stator 3. The experimental results show a dominant casing secondary flow loss core that has migrated to 50% span and a smaller hub loss core at 25% span. With no leakage flows (clean hub and casing annulus lines), TURBOSTREAM predicts two distinct small loss cores at 15% and 85% span. With the addition of leakage paths, the agreement

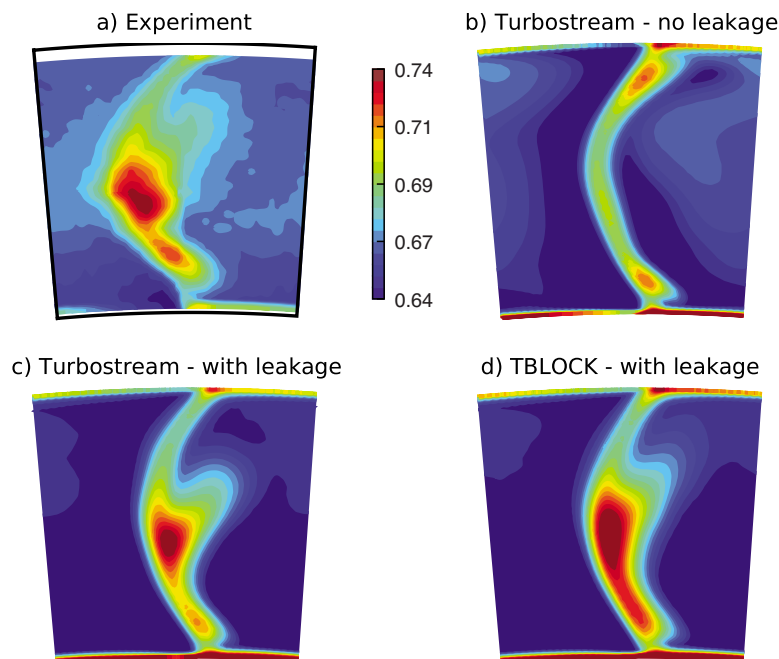


Fig. 8 C_{p_0} contours—stator 3

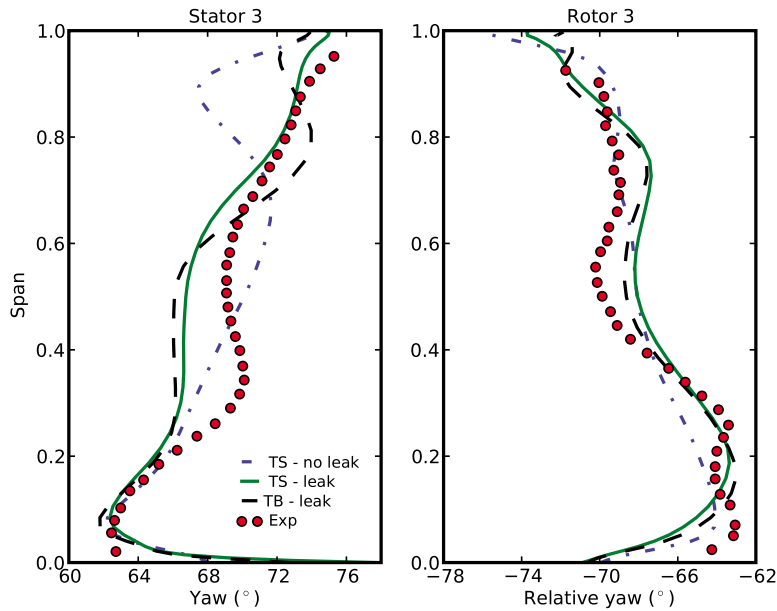


Fig. 9 Measured and predicted pitchwise averaged yaw angle

is much closer. In particular, the shroud leakage from rotor 2 has strengthened the stator 3 casing secondary flow and pushed the associated loss core toward midspan. The remaining discrepancy between the CFD and experiment is likely to be largely the result of difficulties in obtaining the precise leakage gaps in the experiment (particularly at the hub). Finally, Fig. 9 compares the exit yaw angle distributions downstream of stator 3 and rotor 3. Again, the addition of shroud leakage improves the predictions of the stator 3 flow near the casing but, in this case, the accuracy of the rotor exit prediction has not been significantly improved by the inclusion of leakage paths.

7 Discussion

The work presented here has shown that GPUs have enabled a dramatic acceleration of TURBOSTREAM (19 times speed-up on a single GPU versus a quad-core CPU) as compared with the original Fortran solver, TBLOCK. Such a step-change will have two clear implications for the turbomachinery design process:

First, as demonstrated by the three-stage turbine calculation presented in this paper, it is now possible to perform steady-state simulations of whole machines in less than 10 min, even on clusters of moderate size and cost. Furthermore, single blade calculations are approaching interactive time scales on desktop computers with a single GPU.

Second, the performance offered by TURBOSTREAM enables the use of high-fidelity methods in the design process. For example, full annulus unsteady simulations, which are not currently routine for design work, can now be done in calculations that can be left to complete overnight.

8 Conclusions

The main conclusion that can be drawn from this work is that massively parallel architectures such as GPUs can provide an order-of-magnitude greater performance than traditional CPUs for CFD solvers. However, taking advantage of processors such as the GPU requires a complete rewrite of the solver. We argue that the rapidly changing many-core processor landscape means that the use of a source-to-source compiler to decouple the solver's definition from its implementation is crucial. This approach has the added benefit of allowing for the use of complicated optimization strategies that would otherwise make it difficult for CFD developers to recognize the underlying algorithm in the source code.

For turbomachinery design, the dramatic increase in perfor-

mance offered by TURBOSTREAM will open up new levels of interactivity in three-dimensional design, as well as enabling the use of high-fidelity methods in the routine design process.

Acknowledgment

The authors would like to thank NVIDIA for donating the GPU hardware used in this work. In addition, the authors are grateful to Budimir Rosic of the Whittle Laboratory for providing the three-stage turbine test case.

Nomenclature

\mathbf{u}	= velocity vector
λ	= thermal conductivity
ρ	= density
τ	= viscous stress tensor
c_v	= specific heat at constant volume
R	= gas constant
e	= specific entergy = $c_v T + \frac{1}{2} v^2$
h_0	= stagnation enthalpy
p	= pressure
p_0	= stagnation pressure
T	= temperature
t	= time
v	= velocity magnitude
Δs	= entropy change

Appendix: Implementation Details

To illustrate how the implementation of stencil operations differs for CPUs and GPUs, we will consider the simple second-order smoothing stencil defined below.

$$b_{i,j,k} = (1-s)a_{i,j,k} + \frac{s}{6}(a_{i-1,j,k} + a_{i+1,j,k} + a_{i,j-1,k} + a_{i,j+1,k} + a_{i,j,k-1} + a_{i,j,k+1}) \quad (A1)$$

where a and b are values in a structured grid indexed by i, j, k , and s is a factor controlling the amount of smoothing.

To simplify the problem, we only consider a computational domain consisting of a single block. The block has the dimensions NI-2, NJ-2, and NK-2 in the three coordinate directions. In memory, this block is represented as a three-dimensional array

with dimensions NI, NJ, and NK, where the extra two points in each dimension contain ghost cells, one on either side of the domain in each dimension. These ghost cells are assumed to contain the appropriate values so that whatever boundary conditions exist around the block are satisfied when we perform the stencil operation at the edges of the domain.

Listings 1 and 2 contain examples of the implementation of the stencil for a CPU and an NVIDIA GPU, respectively. The CPU implementation is in Fortran 77; the GPU implementation is in NVIDIA's CUDA language. The examples include the memory allocation, the calling of the subroutine and the definition of the subroutine itself (note that subroutines are referred to as kernels on the GPU). For the sake of brevity, we do not show the initialization of the memory. The CPU implementation should be familiar to most CFD developers. It consists of a simple nested loop over the computational domain, with the inner computation performing the stencil operation.

Listing 1 Fortran implementation for a CPU

```

REAL A(NI,NJ,NK), B(NI,NJ,NK), SF
CALL SMOOTH(SF, A, B)

SUBROUTINE SMOOTH(SF, A, B)
C
  LOOP OVER DOMAIN
  DO K=2,NK-1
    DO J=2,NJ-1
      DO I=2,NI-1
        B(I,J,K)=(1.0-SF)*A(I+1,J,K)+
        & SF*(A(I-1,J,K)+A(I+1,J,K)+
        & A(I,J-1,K)+A(I,J+1,K)+
        & A(I,J,K-1)+A(I,J,K+1))/6.0
      END DO
    END DO
  END DO
RETURN
END

```

Listing 2 CUDA implementation for an NVIDIA GPU

```

1  /* Macro for 3D to 1D index translation */
2  #define I3D(ni,nj,i,j,k) ((i)+(ni)*(j)+(ni)*(nj)*(k))
3
4  float *a_h, *a_d, *b_d, sf; /* variable declarations */
5
6  /* allocate memory on the host (CPU) */
7  nbyte = sizeof(float)*NI*NJ*NK;
8  a_h = malloc(nbyte);
9  /* allocate memory on the device (GPU) */
10 cudaMalloc(&a_d, nbyte);
11 cudaMalloc(&b_d, nbyte);
12 /* transfer memory from host to device */
13 cudaMemcpy(a_d, a_h, nbyte, cudaMemcpyHostToDevice);
14
15 /* GPU kernel parameters */
16 num_threadblocks = dim3(1,1,1); /* single thread block */
17 num_threads = dim3(NI,NK,1); /* plane of threads */
18 /* invoke GPU kernel */
19 smooth_kernel<<<num_threadblocks, num_threads>>>(a_d, b_d, sf);
20
21 __global__ void smooth_kernel(float sf, float *a_d, float *b_d)
22 {
23   int i, j, jml, jp1, k, j-plane;
24   __shared__ float a[NI][3][NK]; /* shared memory for three planes */
25
26   i = (int)threadIdx.x; /* current thread index */
27   k = (int)threadIdx.y;
28
29   /* fetch the first planes into shared memory */
30   a[i][0][k] = a_d[I3D(NI, NJ, i, 0, k)];
31   a[i][1][k] = a_d[I3D(NI, NJ, i, 1, k)];
32
33   jml = 0; j = 1; jp1 = 2; /* set initial jml, j, jp1 */
34
35   /* iterate upwards in j-direction */
36   for (j-plane=1; j-plane < NJ-1; j-plane++) {
37
38     /* read the next plane into the jp1 slot */
39     a[i][jp1][k] = a_d[I3D(NI, NJ, i, j-plane+1, k)];
40     /* make sure reads into shared memory are done */
41     __syncthreads();
42
43     /* ghost-zone threads don't compute */
44     if (i > 0 && i < ni-1 && k > 0 && k < nk-1) {
45       /* apply stencil and write out result */
46       i000 = I3D(NI, NJ, i, j, k);
47       b_d[i000] = (1.0f-sf)*a[i][j][k] + sf*(a[i-1][j][k] + a[i+1][j][k] +
48       a[i][jml][k] + a[i][jp1][k] + a[i][j][k-1] + a[i][j][k+1])/6.0f;
49     }
50     tmp = jml; jml = j; j = jp1; jp1 = tmp; /* cycle j indices */
51   }
52 }

```


The GPU implementation is more complicated. First, there are now two disjoint memory spaces to manage, one for the CPU and one for the GPU. It is therefore necessary to allocate memory on both the CPU (line 8) and the GPU (lines 10 and 11) and then transfer data from the CPU to the GPU (line 13). Any operations that involve GPU memory outside of a kernel require calls to special functions implemented by the NVIDIA GPU driver, these have the prefix *cuda*. Second, since a GPU kernel is executed in parallel by many threads at the same time, it is necessary when calling it to specify how many threads are needed and how these should be organized. Here, we assume that the domain is small enough to fit in a single sub-block, so that only a single “thread block” is needed (line 16). A thread block is a CUDA term for a group of threads that operate together, and are executed in parallel on the same multiprocessor. In this particular example, the thread block consists of a single plane of threads (line 17). Note that CUDA’s facility for multidimensional thread blocks is used to simplify the indexing in the kernel later. Finally, the kernel is called with the required number of threads (line 19).

Regarding the implementation of the kernel itself, the following points should be noted.

1. Two CUDA-specific keywords are used. The GPU kernel is defined as `__global__` (line 21), which means that it is called from the CPU and is executed on the GPU. The array storage in the kernel is defined as `__shared__` (line 24), which means that the arrays are stored in the 16 KB on-chip memory associated with each of the GPU’s multiprocessors.
2. Each thread uses the built-in variable `threadIdx` to find its *i* and *k* coordinates in the plane of a sub-block (lines 26 and 27).
3. The variables `jm1`, `j`, and `jp1` are used to hold the offsets to the *j*−1, *j*, and *j*+1 planes in shared memory (line 33). These are cycled at the end of each iteration (line 50) so that the new plane that is loaded during the next iteration replaces the one that is no longer required by the stencil operation.
4. Data from the array `a_d` in the GPU’s main memory is explicitly loaded into the array `a` in shared memory (lines 30, 31, and 39). The threads in the plane load one value each, so the built-in function `__syncthreads()` has to be called to make sure all the threads have finished loading data before progressing further in the code.
5. The outer threads only load data from the ghost zones and do not participate in the computation with the inner threads (line 44).

References

- [1] Brandvik, T., and Pullan, G., 2007, “Acceleration of a Two-Dimensional Euler Solver Using Commodity Graphics Hardware,” *J. Mech. Eng. Sci.*, **221**(12), pp. 1745–1748.
- [2] Brandvik, T., and Pullan, G., 2008, “Acceleration of a 3D Euler Solver Using Commodity Graphics Hardware,” AIAA Paper No. 2008-607.
- [3] Elsen, E., LeGresley, P., and Darve, E., 2008, “Large Calculation of the Flow Over a Hypersonic Vehicle Using a GPU,” *J. Comput. Phys.*, **227**(24), pp. 10148–10161.
- [4] Borkar, S., 2007, “Thousand Core Chips: A Technology Perspective,” *Proceedings of the 44th Annual Conference on Design Automation*.
- [5] Klostermeier, C., 2008, “Investigation Into the Capability of Large Eddy Simulation for Turbomachinery Design,” Ph.D. thesis, University of Cambridge.
- [6] Reid, K., Denton, J. D., Pullan, G., Curtis, E., and Longley, J., 2007, “The Interaction of Turbine Inter-Platform Leakage Flow With the Mainstream Flow,” *ASME J. Turbomach.*, **129**(2), pp. 303–310.
- [7] Rosic, B., Denton, J. D., and Pullan, G., 2006, “The Importance of Shroud Leakage Modeling in Multistage Turbine Calculations,” *ASME J. Turbomach.*, **128**(4), pp. 699–707.
- [8] Denton, J. D., 1975, “A Time Marching Method for Two- and Three-Dimensional Blade to Blade Flows,” Aeronautical Research Council Reports and Memoranda, Report No. 3775.
- [9] Denton, J. D., 1982, “An Improved Time Marching Method for Turbomachinery Flow Calculation,” *ASME Paper No. 82-GT-239*.
- [10] Denton, J. D., 1986, “The Use of a Distributed Body Force to Simulate Viscous Effects in 3D Flow Calculations,” *ASME Paper No. 86-GT-144*.
- [11] Denton, J. D., 1990, “The Calculation of Three Dimensional Viscous Flows Through Multistage Turbines,” *ASME Paper No. 90-GT-19*.
- [12] Denton, J. D., 2002, “The Effects of Lean and Sweep on Transonic Fan Performance,” *TASK Q.*, **6**(1), pp. 7–23.
- [13] Tucker, P. G., Rumsey, C. L., Spalart, P. R., Bartels, R. E., and Biedron, R. T., 2005, “Computations of Wall Distances Based on Differential Equations,” *AIAA J.*, **43**(3), pp. 539–549.
- [14] Denton, J. D., and Singh, U. K., 1979, “Time Marching Methods for Turbomachinery Flow Calculation,” von Karman Institute for Fluid Dynamics, Application of Numerical Methods to Flow Calculations in Turbomachines.
- [15] Dawes, W. N., 1992, “Toward Improved Throughflow Capability: The Use of Three-Dimensional Viscous Flow Solvers in a Multistage Environment,” *ASME J. Turbomach.*, **114**(1), pp. 8–17.
- [16] Fritsch, G., and Giles, M. B., 1992, “Second-Order Effects of Unsteadiness on the Performance of Turbomachines,” *ASME Paper No. 92-GT-389*.
- [17] Jameson, A., 1991, “Time Dependent Calculations Using Multigrid, With Applications to Unsteady Flows Past Airfoils and Wings,” AIAA Paper No. 91-1596.
- [18] Datta, K., Murphy, M., Volkov, V., Williams, S., Carter, J., Oliker, L., Patterson, D., Shalf, J., and Yelick, K., 2008, “Stencil Computation Optimization and Autotuning on State-of-the-Art Multicore Architectures,” *Proceedings of Supercomputing 2008*.
- [19] Williams, S., Shalf, J., Oliker, L., Kamil, S., Husbands, P., and Yelick, K., 2007, “Scientific Computing Kernels on the Cell Processor,” *Int. J. Parallel Program.*, **35**(3), pp. 263–298.