

Acceleration of a CFD code with a GPU

Dennis C. Jespersen

NASA/Ames Research Center, Moffett Field, CA, USA

Tel.: +1 650 604 6742; Fax: +1 650 604 4377; E-mail: Dennis.Jespersen@nasa.gov

Abstract. The Computational Fluid Dynamics code OVERFLOW includes as one of its solver options an algorithm which is a fairly small piece of code but which accounts for a significant portion of the total computational time. This paper studies some of the issues in accelerating this piece of code by using a Graphics Processing Unit (GPU). The algorithm needs to be modified to be suitable for a GPU and attention needs to be given to 64-bit and 32-bit arithmetic. Interestingly, the work done for the GPU produced ideas for accelerating the CPU code and led to significant speedup on the CPU.

Keywords: GPU, CUDA, acceleration, OVERFLOW code

1. Introduction

Computational Fluid Dynamics (CFD) has a history of seeking and requiring ever higher computational performance. This quest has in the past partly been satisfied by faster CPU clock speeds. The era of increasing clock rates has reached a plateau, due mainly to heat dissipation constraints. A boost in computational performance without increasing clock speed can be supplied by parallelism. This parallelism can come in the form of task parallelism, data parallelism or perhaps a combination of the two. Common current paradigms for implementing parallelism are explicit message-passing with MPI [6] for either distributed or shared memory systems and OpenMP [7] for shared memory systems. A hybrid paradigm is also possible, using OpenMP within multiprocessor nodes and MPI for inter-node communication.

A Graphics Processing Unit (GPU) is a processor specialized for graphics rendering. The development of GPUs has been mostly driven by the demand for fast high-quality graphics images for computer games. The large market for computer games provided resources for rapid development of more powerful GPUs. The nature of the graphical display task, highly parallel rendering of many pixels, drove the development of parallel hardware with high computational capability.

GPUs were originally programmed in specialized graphics programming languages. A few researchers realized that some non-graphical computing problems could be expressed as graphical programming tasks, allowing the power of the GPU to be used in a non-graphical environment. Using a GPU in this manner

was difficult, with a steep learning curve. In order to make GPU programming more accessible, several projects were initiated to develop programming languages that could exploit the power of GPUs. Some efforts in this area are Brook and BrookGPU [10] and CUDA [5]. A standardization effort is OpenCL (Open Computing Language) [9].

The enhanced accessibility of GPUs has led to much recent work in the area of general-purpose computing on GPUs. A GPU can produce a very high flop/s (floating-point operations per second) rate if an algorithm is well suited for the device. There have been several studies illustrating the acceleration of scientific computing codes that is possible by using GPUs [2, 13, 17]. In this paper we study the issues in accelerating a well-known CFD code, OVERFLOW, on a GPU.

The viewpoint taken here is that the GPU acts as a co-processor to the CPU. The contemporary CPU typically is a 4-core processor. Serial or modestly parallel (less than 10 threads, say) parts of a code should be executed on the CPU while massively parallel parts of a code should be executed on the GPU. The paradigm is: prepare data on the CPU, transfer to the GPU, execute on the GPU and transfer results back from the GPU to the CPU. Note that this paradigm implies that performance measures and timing comparisons between a pure CPU and a CPU + GPU combination should include the time required for data transfer to and from the GPU.

A characteristic of GPUs inherited from their use strictly as graphics rendering engines is the weak support for 64-bit arithmetic. At the time this work was performed, 64-bit arithmetic was an order of magni-

tude slower than 32-bit arithmetic. This was a strong factor impelling the examination of judicious use of 32-bit arithmetic.

2. GPU environment

For our purposes here the key issues of GPUs are massive parallelism, ideally with thousands of threads, 32-bit floating-point arithmetic, and the overhead of data traffic between the CPU and GPU. For a GPU to successfully accelerate a code segment the code must be amenable to large-scale parallelism, must contain enough computational work to amortize the cost of transferring data from the CPU to the GPU and transferring results back from the GPU and should tolerate 32-bit floating-point arithmetic. (Recent GPU hardware supports some 64-bit arithmetic but 32-bit arithmetic is significantly faster.) We will see that the SSOR algorithm in OVERFLOW is not well suited to a GPU, but that a Jacobi version of the algorithm might be suitable for a GPU.

This work used CUDA (Compute Unified Device Architecture) to support the use of the GPU. CUDA, developed by NVIDIA Corp., is a combination of hardware and software that enables several models of NVIDIA graphics cards to be used as general-purpose processors. The “device” (GPU) is physically a set of multiprocessors, say on the order of 20 multiprocessors, each of which is itself a set of cores, perhaps 8–32 cores, giving a total of a few hundred cores. Logically, the device is organized as a one- or two-dimensional array of blocks and each block is structured as a one-, two- or three-dimensional array of threads. All the threads in a specific block execute in a specific multiprocessor. All threads execute the same code, each thread operating on its own data, so this is a data-parallel programming paradigm. There is a global memory accessible to all threads. In addition, threads in a given multiprocessor can cooperate via a “shared memory” (effectively a cache). Threads in a given multiprocessor can synchronize with one another but there is no synchronization across multiprocessors. The shared memory is small and accesses to it have low latency. The global memory is large but with high latency. One can increase bandwidth to global memory by properly grouping (“coalescing”) memory requests. One hopes that global memory latency can be covered by using a large number of threads and switching quickly from one thread to another if a thread is blocked waiting for data.

The actual programming language for the CUDA architecture is C with a small set of extensions. These extensions consist of keywords to describe where a function executes and where data reside. In addition each thread has predefined variables that describe its location in its batch and the location of its batch in the grid of multiprocessors. Finally, there is syntax denoting the execution of a “kernel” (function which executes on the device) and the configuration of the device for the kernel. So a CUDA program consists of C code interspersed with calls to one or more kernel functions. There is library support for device initialization, transfer of data to and from the device and other miscellaneous activities. An example of a CUDA program is given in the Appendix.

Attractive features of CUDA are the gentle learning curve, wide availability of CUDA-capable devices, and large and active user community. A potential weakness is the tie to a single vendor, resulting in lack of portability. The OpenCL project [9] is an open standard for parallel programming of heterogeneous systems. When this work was begun the OpenCL standard had not been finalized.

A significant feature of CUDA programming, which might be seen as either a strength or weakness, is the necessity for the programmer to specify the placement of data in global and shared memory, which effectively amounts to a requirement to explicitly manage the cache. This goes along with the necessity to explicitly define the layout of the GPU (number of blocks, number of threads) and gives CUDA programming a distinct feeling of writing at a low level, close to the hardware. The mapping of a numerical algorithm onto the GPU can often be defined in a variety of ways, each of which might give different performance, and the only way to evaluate the performance of the various alternatives is by actual coding and testing.

3. OVERFLOW code

The OVERFLOW code [4,11,12,16] is intended for the solution of the Reynolds-averaged Navier–Stokes equations with complex geometries. The code uses finite differences on logically Cartesian meshes. The meshes are body-fitted and geometric complexity is handled by allowing the meshes to arbitrarily overlap one another.

OVERFLOW uses implicit time-stepping and can be run in time-accurate or steady-state mode. Implicit time-stepping is used because implicit methods tend to

mitigate severe stability limits on the size of the time step that arise for explicit time-stepping methods on highly-stretched grids; such grids are common for viscous flow problems at high Reynolds numbers. A consequence of implicit time-stepping is that some method is needed to (approximately) solve the large system of equations that arises when advancing from one time level to the next.

The OVERFLOW user needs to specify physical flow inputs, such as Mach number and Reynolds number, and boundary conditions which typically define solid walls and inflow or outflow regions. Along with these physics-type inputs there are inputs which choose particular numerical algorithms and specify parameters for them.

The basic equation of fluid motion solved by OVERFLOW is of the form:

$$\partial Q/\partial t + L(Q) = f(Q), \quad (1)$$

where Q is the vector of flow variables, $L(Q)$ denotes all the spatial differencing terms, and $f(Q)$ denotes terms from boundary conditions and possible source terms.

Equation (1) is discretized and written in “delta form” [1,15]:

$$A(\Delta Q^{n+1}) = R(Q^n), \quad (2)$$

where $A = A(Q^n)$ is a very large sparse matrix which is not explicitly constructed, $\Delta Q^{n+1} = Q^{n+1} - Q^n$ is the vector of unknowns, and $R(Q^n)$ involves the discretization of the $L(Q)$ and $f(Q)$ terms at time level n . The user of OVERFLOW must choose among several possible discretizations (e.g., central differencing, Total Variation Diminishing, Roe upwind). Each of these choices typically requires further user specification of numerical parameters such as dissipation parameters or type of flux limiter and parameters for the limiter. Finally, the user needs to decide which implicit algorithm to use: some choices are approximate factored block tridiagonal, approximate factored scalar pentadiagonal or LU-SGS (approximate LU factorization with Symmetric Gauss–Seidel iteration). Over the years the code evolved and expanded to incorporate six basic choices for the implicit part of the algorithm.

4. The SSOR algorithm in OVERFLOW

In an attempt to ease the user’s burdensome task of selecting algorithm options and choosing parameters

which may change for each class of flow problem, recently another option was added for the implicit part of OVERFLOW, with the hope that it would be widely applicable and would be almost universally usable [14]. In the reference, this algorithm is referred to as an SSOR algorithm, but it is strictly speaking a mix of an SSOR algorithm and a Jacobi algorithm, so it might be called a “quasi-SSOR” algorithm.

The key step of the algorithm is as follows. At each grid point with index (j, k, l) one computes a residual R_{jkl} and six 5×5 matrices A^+ , B^+ , C^+ , A^- , B^- , C^- ; these matrices depend on the flow variables at the neighboring grid points and are fixed during the SSOR iterations. Then, with iteration stage denoted by a superscript n and with a relaxation parameter ω , relaxation steps are of the form:

$$\begin{aligned} \Delta Q_{jkl}^{n+1} = & (1 - \omega)\Delta Q_{jkl}^n \\ & + \omega(R_{jkl} - A_{jkl}^+ \Delta Q_{j-1,k,l}^n \\ & - B_{jkl}^+ \Delta Q_{j,k-1,l}^{n+1} - C_{jkl}^+ \Delta Q_{j,k,l-1}^{n+1} \\ & - A_{jkl}^- \Delta Q_{j+1,k,l}^n - B_{jkl}^- \Delta Q_{j,k+1,l}^n \\ & - C_{jkl}^- \Delta Q_{j,k,l+1}^n) \end{aligned} \quad (3)$$

for a forward sweep (assuming the 5-vectors $\Delta Q_{j,k-1,l}^{n+1}$ and $\Delta Q_{j,k,l-1}^{n+1}$ have been computed, and updating all ΔQ_{jkl} as soon as a full line of j values has been computed) and a step of the form:

$$\begin{aligned} \Delta Q_{jkl}^{n+1} = & (1 - \omega)\Delta Q_{jkl}^n \\ & + \omega(R_{jkl} - A_{jkl}^+ \Delta Q_{j-1,k,l}^n \\ & - B_{jkl}^+ \Delta Q_{j,k-1,l}^n - C_{jkl}^+ \Delta Q_{j,k,l-1}^n \\ & - A_{jkl}^- \Delta Q_{j+1,k,l}^n - B_{jkl}^- \Delta Q_{j,k+1,l}^{n+1} \\ & - C_{jkl}^- \Delta Q_{j,k,l+1}^{n+1}) \end{aligned} \quad (4)$$

for a backward sweep (again assuming $\Delta Q_{j,k+1,l}^{n+1}$ and $\Delta Q_{j,k,l+1}^{n+1}$ have been computed). The forward/backward pair is then iterated. This algorithm is not strictly speaking an SSOR algorithm; it is Jacobi in j and SSOR in k and l . We will refer to it as SSOR for simplicity. This algorithm needs the values of ΔQ at the six nearest spatial neighbors of grid point (j, k, l) , some of them at iteration level n and some of them at iteration level $n + 1$.

The SSOR algorithm is a modest-sized subroutine but testing shows it may consume 80% of the total runtime of the code, so it is a computational hot spot. The modest size of the subroutine and the large fraction of total time consumed by the algorithm make using a GPU as a co-processor to accelerate the code an attractive idea.

Unfortunately, the algorithm as it stands is not suited to a GPU due to the dependencies of the iteration, namely ΔQ^{n+1} appears on the right-hand side of Eqs (3) and (4). An algorithm that would be suited to a GPU would be a Jacobi algorithm with relaxation steps of the form:

$$\begin{aligned} \Delta Q_{jkl}^{n+1} = & (1 - \omega)\Delta Q_{jkl}^n \\ & + \omega(R_{jkl} - A_{jkl}^+ \Delta Q_{j-1,k,l}^n \\ & - B_{jkl}^+ \Delta Q_{j,k-1,l}^n - C_{jkl}^+ \Delta Q_{j,k,l-1}^n \\ & - A_{jkl}^- \Delta Q_{j+1,k,l}^n - B_{jkl}^- \Delta Q_{j,k+1,l}^n \\ & - C_{jkl}^- \Delta Q_{j,k,l+1}^n). \end{aligned} \quad (5)$$

Here we could envision assigning a thread of computation to each grid point and the threads could com-

pute independently of one another because there are no ΔQ^{n+1} terms on the right-hand side of (5).

It is important to realize that the Jacobi algorithm might be less robust or might converge slower than the original SSOR algorithm. Fully discussing this would take us too far afield, though we will show some convergence comparisons of Jacobi and SSOR.

The work presented here proceeded in several stages:

1. Implement a Jacobi algorithm on the CPU using 64-bit arithmetic; compare performance and convergence/stability of Jacobi and SSOR.
2. Implement a Jacobi algorithm on the CPU using 32-bit arithmetic; compare performance and convergence/stability of 64-bit and 32-bit Jacobi.
3. Implement a Jacobi algorithm on the GPU; compare performance of the GPU algorithm with the 32-bit CPU algorithm.

5. Implementation and results

The first stage of the work, implementing the Jacobi algorithm on the CPU using 64-bit arithmetic, was straightforward. We compare in Fig. 1 convergence for the SSOR and Jacobi algorithms on two test cases. The

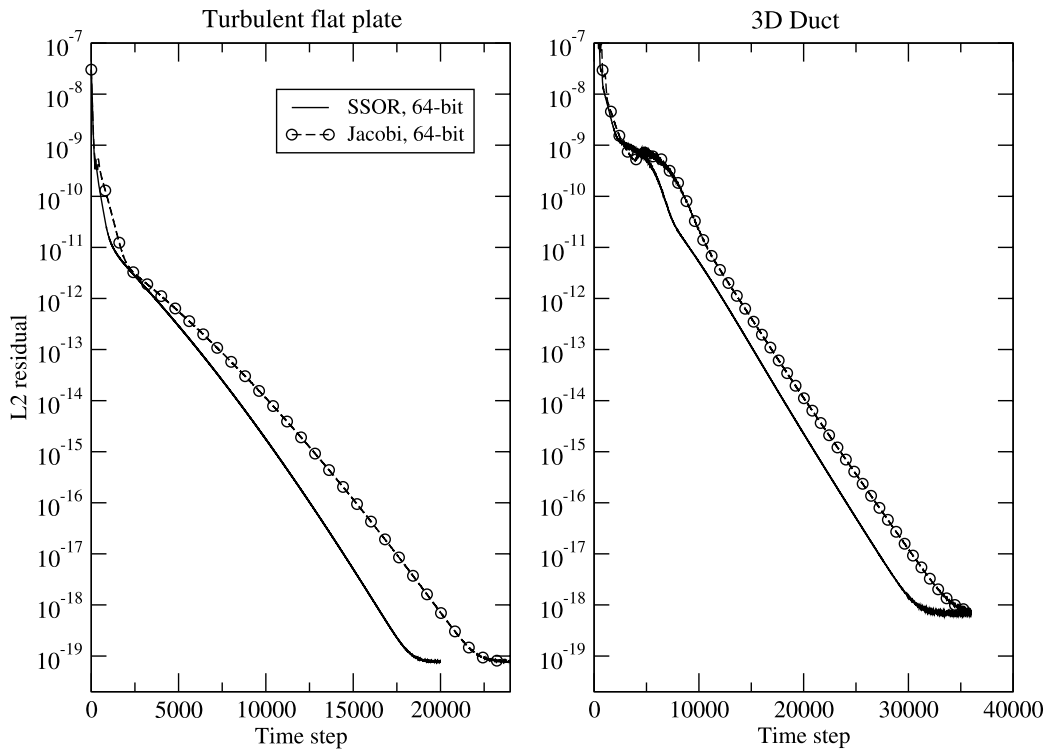


Fig. 1. SSOR and Jacobi convergence, 64-bit arithmetic.

first test case is turbulent flow over a flat plate with a $121 \times 41 \times 81$ grid. The second flow is turbulent flow in a curved duct with a $166 \times 31 \times 49$ grid. Both cases show, unsurprisingly, that asymptotic convergence of the Jacobi algorithm is slightly slower than that of the SSOR algorithm. Both cases reach machine zero (solution converged to 64-bit accuracy). The SSOR algorithm is slightly faster in terms of wallclock seconds per time step, because the SSOR algorithm updates ΔQ as the computation proceeds whereas the Jacobi algorithm uses an extra array to store the changes to ΔQ and then sweeps through the full ΔQ array to form the new values of ΔQ .

Implementing the Jacobi algorithm in 32-bit arithmetic for the CPU was tedious but straightforward. The implementation included making 32-bit versions of all the subroutines dealing with the computation of the left-hand side matrices (about 50 subroutines) and copying, on the front-end, the flow variables and metric terms to 32-bit quantities; in all, 28 words per grid point were copied from 64-bit to 32-bit representation. In Fig. 2 we show convergence for the Jacobi algorithm in 64-bit arithmetic and in 32-bit arithmetic for the two test cases. To plotting accuracy there is no difference in convergence between the 32-bit and 64-bit Jacobi al-

gorithms. This verifies for these cases that full 64-bit solution accuracy can be obtained with a 32-bit Jacobi algorithm.

Finally, the Jacobi algorithm was implemented on the GPU. The strategy was to compute all the matrices A^+ , etc., on the CPU and transfer them to the GPU. The Jacobi algorithm itself, just one subroutine, was hand-translated into CUDA code. This strategy avoided a long error-prone translation of many Fortran subroutines into CUDA, but this strategy may be suboptimal as the matrices themselves could be computed on the GPU. We found no difference between convergence of the 32-bit Jacobi algorithm on the CPU and on the GPU, so the slight differences in details of floating-point arithmetic between the CPU and the GPU have no impact for these cases.

Now we consider performance of the code. The metric we use is wallclock seconds per step, so smaller is better. The GPU algorithm was coded in several slightly different ways, varying in the data layout on the GPU and whether or not shared memory on the GPU was used. Data shown are for the best-performing GPU variant.

The work here was done on two platforms. The first platform was a workstation equipped with a 2.1 GHz

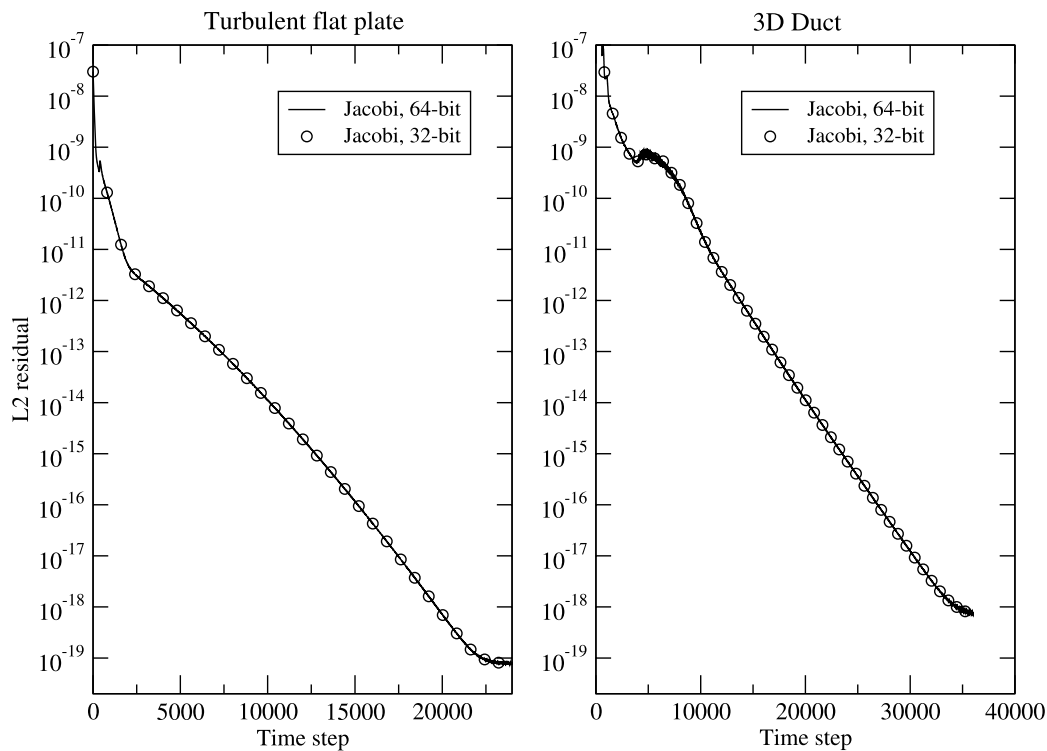


Fig. 2. Jacobi convergence, 64-bit and 32-bit arithmetic.

Table 1
Implicit solver times (lower is better)

Algorithm	G machine		T machine	
	Plate	Duct	Plate	Duct
SSOR CPU (s/step)	3.51	2.14	3.83	2.33
Jacobi GPU (s/step)	1.43	0.91	1.35	0.76
GPU/CPU ratio	0.41	0.43	0.35	0.33

Table 2
Total time for CPU and GPU (lower is better)

Algorithm	G machine		T machine	
	Plate	Duct	Plate	Duct
SSOR CPU (s/step)	6.96	4.21	7.93	4.85
Jacobi GPU (s/step)	4.41	2.66	5.04	3.12
GPU/CPU ratio	0.63	0.63	0.64	0.64

quad-core AMD Opteron 2352 processor. The host compiler system was the Portland Group compiler suite version 8. The GPU card was a 1.35 GHz NVIDIA GeForce 8800 GTX with 128 cores and 768 MB of global memory. The connection between CPU and GPU was a PCI Express 16X bus. The programming interface was CUDA version 1.0. This platform will be referred to as the “G machine”.

The second platform was a workstation equipped with two 2.8 GHz dual-core AMD Opteron 2220 processors. The GPU card was a 1.30 GHz NVIDIA Tesla C1060 with 240 cores and 4 GB of global memory. For this machine, the source code was cross-compiled on the first machine using the Portland Group compiler. This platform will be referred to as the “T machine”.

Tables 1 and 2 give performance data for the two test cases on the two machines. The implicit solver times in Table 1 (which include a small amount of work on the CPU as well as the actual relaxation algorithm) show a speedup on the GPU by about a factor of between 2.5 and 3. The reason the T machine times are only slightly better than the G machine times is that the times shown here include some CPU work and for some unknown reason the CPU routines involved ran faster on a single CPU of the G machine than on a single CPU of the T machine. The wallclock time for the full code, which is the quantity of ultimate interest to the code user, decreases by about 40%, as seen in Table 2. Again, a single CPU of the G machine is overall faster than a single CPU of the T machine.

Table 3 gives GPU total time (kernel plus time for data transfer) and GPU kernel time for these cases. For the 8800 GTX device, the kernel which gave the best

Table 3
GPU kernel and data transfer times, s/step (lower is better)

	8800 GTX		Tesla C1060	
	Plate	Duct	Plate	Duct
GPU total	0.904	0.576	0.314	0.193
GPU kernel only	0.784	0.499	0.142	0.082
Data transfer	0.120	0.077	0.172	0.111

overall code performance was a kernel which mapped each grid point to a different thread on the GPU (thanks to Jonathan Cohen of NVIDIA for showing a nice way to do this) and which used some shared memory. For the Tesla device, the kernel which gave the best overall code performance was a kernel involving a two-dimensional mapping of the first two grid dimensions onto the device, a loop in the 3rd dimension and 16 threads per grid point with the 5×5 matrices on the CPU stored in an array of size 32. For both the Tesla and the 8800 GTX devices there are data layouts which give better performance of the GPU considered in isolation, but these layouts involve data motion on the CPU and this data motion loses more wallclock time than is gained by the faster kernel. From this table it can be seen that the time for data transfer to and from the GPU is relatively small for the 8800 GTX but is significant (more than the time for the kernel itself) for the C1060.

6. Impact of GPU work on CPU code

These results are encouraging. It seems that the Jacobi GPU algorithm is significantly faster than the SSOR CPU algorithm, since Table 2 shows a speedup for the whole code of about 40%. This is a significant speedup considering that the only code being executed on the GPU is a small piece of the implicit side and there are no changes to the flow explicit side or to the turbulence model.

However, further reflection indicates that more performance can be gained from the CPU. Specifically, the SSOR CPU algorithm could be changed to use 32-bit arithmetic. Experience has shown that OVERFLOW is a bandwidth-limited code, so use of 32-bit arithmetic (where applicable) should significantly speed up the code due to the effective doubling of bandwidth. This motivated the implementation of a 32-bit SSOR algorithm on the CPU. Convergence for this algorithm, for the plate and duct test cases,

Table 4

SSOR performance on CPU, 64-bit and 32-bit (lower is better)

Algorithm	G machine		T machine	
	Plate	Duct	Plate	Duct
SSOR-64 CPU (s/step)	6.96	4.21	7.93	4.85
SSOR-32 CPU (s/step)	5.55	3.34	6.30	3.87
SSOR-32/SSOR-64 ratio	0.80	0.79	0.79	0.80

Table 5

SSOR and OpenMP performance on CPU, s/step (lower is better)

Algorithm	OpenMP threads	G machine		T machine	
		Plate	Duct	Plate	Duct
SSOR-64	1	6.96	4.21	7.93	4.85
SSOR-64	2	5.53	3.27	6.76	4.11
SSOR-64	4	4.60	2.80	5.94	3.60
Revised SSOR-64	1	7.79	4.70	8.41	5.14
Revised SSOR-64	2	4.79	2.85	4.76	2.96
Revised SSOR-64	4	3.36	2.04	3.27	1.99
SSOR-32	1	5.55	3.34	6.30	3.87
SSOR-32	2	4.23	2.55	4.72	2.89
SSOR-32	4	3.45	2.11	3.92	2.41
Revised SSOR-32	1	6.07	3.65	6.92	4.20
Revised SSOR-32	2	3.63	2.18	3.79	2.37
Revised SSOR-32	4	2.38	1.47	2.50	1.52

was identical to convergence for the 64-bit SSOR algorithm. Single-core performance for the 32-bit SSOR algorithm is compared with performance for the 64-bit SSOR algorithm in Table 4. Total wallclock time is reduced by 20% by the use of 32-bit arithmetic; this is a significant speedup for such an unsophisticated code modification.

Even more performance can be gained on the CPU by taking advantage of multiple CPU cores. OVERFLOW has long had OpenMP capability, thereby using multiple CPU cores, but the SSOR algorithm as originally coded in OVERFLOW was not amenable to OpenMP parallelism (as mentioned in Section 4, the original coding was Jacobi in the j index and Gauss–Seidel in the k and l indices, while the OpenMP parallelism in OVERFLOW is parallelism in l). Revising the algorithm to be Jacobi in l and Gauss–Seidel in j and k allowed the use of OpenMP for the SSOR algorithm. The revised SSOR algorithm had the same convergence characteristics as the Jacobi algorithm for the duct and plate test cases. The revised algorithm can also be coded in 64-bit or 32-bit arithmetic. We show in Table 5 performance for these algorithms and var-

Table 6

Jacobi on GPU + OpenMP on CPU, performance (lower is better)

No. of OpenMP threads on CPU	8800 GTX		Tesla C1060	
	Plate (s/step)	Duct (s/step)	Plate (s/step)	Duct (s/step)
1	4.39	2.66	4.58	2.85
2	3.10	1.78	2.92	1.83
4	2.32	1.42	1.90	1.18

ious numbers of OpenMP threads; these are all CPU performance numbers, the GPU is not involved here, and this is performance for the full code.

For a single OpenMP thread the revised SSOR algorithm is slower than the original, due to poorer cache utilization, but for 2 or 4 OpenMP threads the revised algorithm is faster than the original.

Finally, Table 6 gives performance data for the two platforms using the Jacobi algorithm on the GPU and OpenMP parallelism on the CPU. To compare the code with and without GPU but otherwise using all the computational resources, the best numbers in Table 5 should be compared with the best numbers in Table 6. The result is that for the workstation with the 8800 GTX GPU, the best time with GPU is just a few percent faster than the best time without GPU, whereas for the workstation with the Tesla C1060 GPU, the best time with GPU is about 25% better than the best time without GPU. The ultimate reason for the better performance of the Tesla C1060 on this code is the relaxed alignment restrictions for coalesced loads as compared to the 8800 GTX.

Even this is not the end of the story, as there are further opportunities for moving computation from the CPU to the GPU. For example, the matrices can be computed in parallel, so this part of the computation, which is now executed by the CPU, can be moved to the GPU. These and other optimizations are currently under investigation.

7. Conclusions

The work presented in this paper has shown a speedup by a factor of between 2.5 and 3 for the SSOR solver in OVERFLOW, and a total wallclock time decrease of about 40%, for a GPU as compared to a single CPU. The GPU work gave ideas and motivation for accelerating the code on multi-core CPUs, so that currently the CPU + GPU code is about 25% faster than the multi-core CPU code.

This study has until now focused on obtaining improved performance with a one CPU + one GPU com-

bination. This is the first step to enhancing OVERFLOW performance via GPUs on realistic problems. However, for almost all realistic cases, OVERFLOW is used with MPI (Message-Passing Interface) and many CPUs. The work here extends naturally to any cluster with a number of multi-core nodes, each node also containing a GPU. OVERFLOW could be used in hybrid mode, with each node corresponding to an MPI process, and each MPI process would have multiple OpenMP threads. The hyperwall-2 at NASA/Ames Research Center [8] is such a cluster and the version of OVERFLOW with GPU capability has run on the hyperwall-2 as a proof of concept.

It is worthwhile to note that the work done here has affected the OVERFLOW code. Official release version 2.1ac of OVERFLOW completely abandoned the 64-bit version of the SSOR algorithm in favor of the 32-bit version, and the revised SSOR algorithm (32-bit arithmetic only) is available as an option. The speedups due to 32-bit arithmetic were so compelling that 64-bit arithmetic is no longer even an option in these portions of the code. This may give food for thought when considering the need for 64-bit arithmetic on GPUs.

Acknowledgements

This work was partially supported by the Fundamental Aeronautics Program of NASA's Aeronautics Research Mission Directorate. NVIDIA Corporation donated the Tesla C1060 hardware.

Appendix

As an example of a CUDA program, consider the simple vector addition $z = x + y$ where x , y and z are vectors of dimension n . In standard C, this might be written as:

```
void s_add_vecs(int n, float *x,
               float *y, float *z) {
    int i;
    for (i=0; i<n; i++)
        z[i] = x[i] + y[i];
}
```

and called via

```
void serial_add_vecs() {
    // define scalar n; allocate and
    // initialize arrays x, y, z
    s_add_vecs(n, x, y, z);
}
```

This could be written in CUDA as follows:

```
__global__
void p_add_vecs(int n, float *x,
               float *y, float *z) {
    int myi = blockIdx.x*blockDim.x
            + threadIdx.x;
    if (myi < n) z[myi] = x[myi] + y[myi];
}
```

with calling code

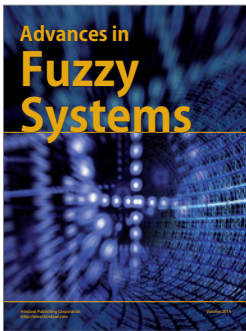
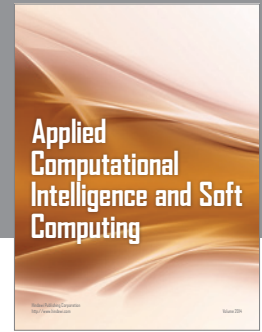
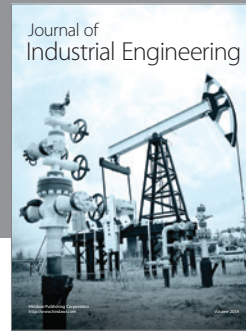
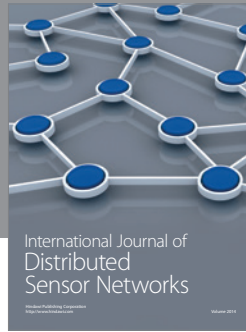
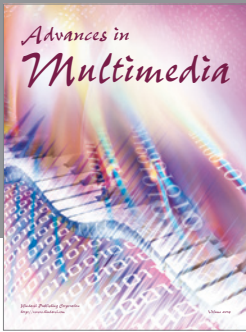
```
void parallel_add_vecs {
    // define scalar n; allocate and
    // initialize arrays x, y, z
    // transfer arrays x, y, z
    // from CPU to GPU
    int nThreads = 256;
    int nBlocks = ceil(n/nThreads);
    p_add_vecs <<<nBlocks,nThreads>>>
        (n, x, y, z);
}
```

The arguments between the triple angle brackets define the execution configuration of the function call. The first argument is a structure specifying the layout of the grid of blocks; in this simple case an integer is silently promoted to a structure specifying a one-dimensional layout. The second argument is a structure specifying the thread layout within each block; again an integer is a shortcut for a one-dimensional layout. So if $n = 100,000$ there would be 391 blocks of 256 threads each. The `__global__` keyword defines `p_add_vecs` as a function which will execute on the GPU. The location of the variables x , y , z is not specified, so by default they reside in the global memory of the GPU. The structures `blockIdx`, `blockDim` and `threadIdx` are automatically defined when the kernel function is called and serve to define the logical organization of the device. Given these structures, each thread computes an index `myi` into the global arrays, and if that index is less than the dimension of the arrays, the thread does work. With $n = 100,000$ there are 100,096 threads and the last 96 threads do no work, but this is much less than 1% of the total number of threads. The number 256 of threads here is just one possibility. Each particular model of NVIDIA card has defined upper limits on the grid and block dimensions and the number of threads per block. In addition, each thread requires some hardware resources, so for any given kernel there may be further restrictions. Meanwhile, each block should have “many” threads for effective parallel use of the available resources. Optimal choice of the number of threads typically requires some experimentation for each particular problem. In our CFD application, the best choice for number of

threads and blocks might depend on the dimensions of the CFD spatial grid.

References

- [1] R. Beam and R.F. Warming, An implicit finite-difference algorithm for hyperbolic systems in conservation law form, *J. Comp. Physics* **22**(1) (1976), 87–110.
- [2] T. Brandvik and G. Pullan, Acceleration of a 3D Euler solver using commodity graphics hardware, in: *AIAA 46th Aerospace Sciences Meeting*, Reno, NV, USA, 2008, paper no. AIAA-2008-607.
- [3] I. Buck, T. Foley, D. Horn, J. Suferman, K. Fatahalian, M. Houston and P. Hanrahan, Brook for GPUs: stream computing on graphics hardware, in: *SIGGRAPH 2004*, 2004.
- [4] P.G. Buning, I.T. Chiu, S. Obayashi, Y.M. Rizk and J.L. Steger, Numerical simulation of the integrated space shuttle vehicle in ascent, in: *AIAA Atmospheric Flight Mechanics Conference*, 1988, paper no. 88-4359-CP.
- [5] http://www.nvidia.com/object/cuda_home.html.
- [6] <http://www-unix.mcs.anl.gov/mpi>.
- [7] <http://www.openmp.org>.
- [8] <http://www.nas.nasa.gov/News/Releases/2008/06-25-08.html>.
- [9] <http://www.khronos.org/opencv>.
- [10] <http://graphics.stanford.edu/projects/brookgpu>.
- [11] W. Kandula and P.G. Buning, Implementation of LU-SGS algorithm and roe upwinding scheme in OVERFLOW thin-layer Navier–Stokes code, in: *AIAA 25th Fluid Dynamics Conference*, Colorado Springs, CO, USA, 1994, paper no. AIAA-94-2357.
- [12] R.L. Meakin, Automatic off-body grid generation for domains of arbitrary size, in: *AIAA 15th Computational Fluid Dynamics Conference*, Anaheim, CA, USA, 2001, paper no. AIAA-2001-2536.
- [13] J. Michalakes and M. Vachharajani, GPU acceleration of numerical weather prediction, *Parallel Processing Letters* **18**(4) (2008), 531–548.
- [14] R.H. Nichols, R.W. Tramel and P.G. Buning, Solver and turbulence model upgrades to OVERFLOW 2 for unsteady and high-speed applications, in: *AIAA 24th Applied Aerodynamics Conference*, San Francisco, CA, USA, 2006, paper no. AIAA-2006-2824.
- [15] T.H. Pulliam and D.S. Chaussee, A diagonalized form of an implicit approximate factorization algorithm, *J. Comp. Phys.* **39**(2) (1981), 347–363.
- [16] K.J. Renze, P.G. Buning and R.G. Rajagopalan, A comparative study of turbulence models for overset grids, in: *AIAA 30th Aerospace Sciences Meeting*, Reno, NV, USA, 1992, paper no. AIAA-92-0437.
- [17] J.C. Thibault and I. Senocak, CUDA implementation of a Navier–Stokes solver on multi-GPU desktop platforms for incompressible flows, in: *AIAA 47th Aerospace Sciences Meeting*, Orlando, FL, USA, 2009, paper no. AIAA-2009-758.



Hindawi

Submit your manuscripts at
<http://www.hindawi.com>

