

# Cell-based Adaptive Mesh Refinement on Hybrid Architectures

David Nicholaeff<sup>\*†§</sup>, Neal Davis<sup>‡</sup>, and Robert Robey<sup>†</sup>

<sup>\*</sup>Oxford University, Oxford, England, United Kingdom

<sup>†</sup>XCP-2 Eulerian Codes, Los Alamos National Laboratory, Los Alamos, NM

<sup>‡</sup>Dept. of Nuclear, Plasma, & Radiological Engineering, University of Illinois at Urbana-Champaign, Urbana, IL

<sup>§</sup>Author emails: David Nicholaeff, dnic@lanl.gov, Neal Davis, davis68@illinois.edu, Robert Robey, brobey@lanl.gov

**Abstract**—Presented in this paper is an OpenCL implementation of a cell-based adaptive mesh refinement (AMR) scheme modeling the two-dimensional shallow water equations using general-purpose graphics processing units (GPGPUs). The challenges associated with ensuring locality of computation in order to fully exploit the throughput and massive data parallelism of the GPU are discussed along with solutions. In particular, data ordering is taken as a free parameter. A stencil-based space-filling curve method allows for optimal load-balancing, while the resulting nontrivial arrangement of cells is addressed by a Cartesian-indexed hash mapping which allows for efficient parallel neighbor accesses. This in turn presents a homogeneous interface to the data across the multiple levels present in heterogeneous architectures.

The relative speed-up of the GPU-enabled AMR code is compared to the respective serial implementation, both with and without the message-passing interface (MPI), providing evidence for the need to design and implement numerical methods on heterogeneous architectures exploiting GPGPUs.

**Keywords**—Cell-based Adaptive Mesh Refinement; AMR; GPGPU; GPU; OpenCL; heterogeneous architecture

Los Alamos National Laboratory is operated by Los Alamos National Security, LLC, for the National Nuclear Security Administration of the U.S. Department of Energy under contract DE-AC52-06NA25396

## I. INTRODUCTION

We present an OpenCL implementation of a cell-based adaptive mesh refinement (AMR) scheme in order to exploit the advantages of the massively parallel general-purpose graphics processing unit (GPGPU). This AMR framework is extensible to a variety of governing equations; we here model a cylindrical dam break problem with the two-dimensional shallow water equations. This representative test case illustrates several of the advantageous capabilities of cell-based adaptive mesh schemes, such as clean symmetry preservation and the treatment of shock-like initial conditions.

Judicious partitioning of the computational domain allows more effective exploitation of the massive number of parallel threads, motivating the structure of algorithm architecture around the goal of achieving maximum data locality. The resulting approach and program are together a proof-of-concept that a cell-based AMR code can be effectively implemented in the memory and threading model provided by OpenCL, which (as of OpenCL 1.1 [1]) provides a device-agnostic

programming language but lacks certain conveniences such as dynamic memory allocation. This last feature, necessary for a properly efficient AMR scheme, is effectuated by a combination of CPU-based memory management and GPU-based computation. Load-balancing is naturally achieved through use of a stencil-based refinement method for array indexing, thereby eliminating the need for a complete recalculation of the mesh array indexing. Finally, a Cartesian-indexed hash mapping scheme to allow fast parallel neighbor accesses at  $O(n)$  is discussed, superseding the use of a  $k$ -D tree at  $O(n \log n)$  (which itself supersedes  $O(n^2)$  algorithms). Comparison of the GPU-enabled AMR code and the original CPU-only serial version, implemented both with and without the message-passing interface (MPI), shows an order of magnitude improvement for GPU-based code.

The results of our cell-based AMR scheme, aptly named CLAMR, provide evidence that parallelization using the GPU delivers significant speed-up for typical numerical simulations and is feasible for scientific applications in the next generation of high-performance computing (HPC). GPU-based computing hides intranode memory latency with increased throughput and hence provides one of the most encouraging paths to exascale computing [2]. The techniques implemented in CLAMR highlight the need for a paradigm shift in programming methodologies that are mindful of parallelism, as analyzed in Davis *et al.* [3]; locality is one such important consideration, and its proper application demands careful reflection on problem decomposition.

This paper is partitioned into five sections. Following the introduction, Section II analyzes the motivation for the architecture decisions behind CLAMR with a review of the adaptive mesh scheme, both logistically and heuristically, as well as a discussion of heterogeneous platforms. The implementation is presented in Section III. This consists of an overview of the numerical method, an overview of the physical model for the shallow water equations, and an analysis of the major challenges of implementation for the adaptive mesh on the GPU – including the difficulties of ensuring proper partitioning of work load and of ensuring locality of memory accesses arising from neighbor searches. Section IV presents our results. The timings between explicit CPU use and a hybrid implementation of the CPU and GPU are compared, both with a single compute node and across compute nodes with the

use of MPI. Ultimately our results suggest that while future numerical physics codes need heterogeneous architectures to see significant speed-up, they cannot do so without devising better algorithms in lieu of ever greater computational power.

## II. CELL-BASED ADAPTIVE MESH REFINEMENT ON HETEROGENEOUS PLATFORMS

### A. The Cell-based AMR Scheme

Numerical models are heavily influenced by their choice of discretization; for example, symmetry preservation is one important feature which can be adversely effected by the structure of the mesh. Hence, a careful selection of the mesh type is critical to properly approximate a continuous space and produce physically legitimate results. We implement a cell-based mesh, a discretization of space into a grid of square cells. The continuous functions of interest, such as mass and momentum, are all taken to be at the center of a cell, which is to say that in the discretization the state variables are averaged over the cell and assigned a Cartesian coordinate existing at the cell center. When the mesh is described as adaptive, it is to say that cells across the mesh can have variable levels of refinement (size) predicated by certain rules.

The motivation for an adaptive mesh is two-fold. First, regions of physical interest are superimposed onto an area of the mesh where the refinement is higher, thus allowing for higher precision. Further, this allows different physical scales to be simultaneously analyzed (e.g. in wave phenomenon long wavelength regions can be placed on coarser cells whereas high frequency waves require greater resolution to discern individual wave peaks)<sup>1</sup>. The second advantage, which directly returns to the motivation to reduce power consumption, is memory frugality. Simply put, the physical model will require far less memory as it will use far less cells for the discretization.

To put the cell-based scheme into perspective, we discuss the types of AMR. The most common AMR scheme in the literature is a structured AMR which superimposes blocks or patches of cells with a finer regular structure over regions of greater physical interest. The technique is described in a series of articles by Berger-Colella-Oliger ([5] and [6]). Perhaps the greatest advantage of this method is that the regular grids in refined regions are just treated like another mesh, making some parts of the implementation much the same as the regular grid. However, it induces additional refinement on cells which could have been left unrefined, thereby not efficiently fulfilling a primary goal of AMR: memory frugality.

The cell-based scheme, on the other hand, refines individual cells, thereby maximizing memory efficiency. As a further consequence, it can precisely refine the regions of the mesh near important physical processes such as shocks or steep wavefronts. And lastly, it has less mesh refinement imprinting for curved or spherical shocks where the regular refined grid

in structured AMR can impact the spherical symmetry of a problem.

On the implementation side of the cell-based AMR scheme, there are several rules dictating the adaptive refinement of the mesh. First, two neighboring cells must have no more than one level of refinement difference. This stems from both ease of implementation and reducing the small error that occurs at refinement steps (the frequency content of the simulated waves will be reduced in half, causing a minor reflection of the wave at the interface). Second, a cell is refined symmetrically, which is to say that it is bisected along all axes. Third, regions of physical interest are refined – this equates to steep gradients in both pressure and material interfaces. Fourth, refinement leads the event, which is to say that refinement should precede before the regions of physical interest arrive. Fifth, indexing is done using a standard Cartesian grid.

### B. Considerations of Heterogeneous Platforms

From a data structures and algorithms perspective, cell-based adaptive mesh refinement for supercomputing applications requires software development methods which consider the dynamics of heterogeneous platforms. For both an overview and references to multiple sources of this argument, see Davis *et al.* [3] The model platform presented here is a simple configuration in which a single CPU communicates with a single GPU, thus defining a single node. These nodes then communicate using MPI. As the numerical calculations are performed on the GPU, we consider this primarily a GPGPU platform. The dynamics of such a platform dictate new considerations in algorithm architecture, which ultimately motivate the design behind the architecture of CLAMR. Production high-performance computing (HPC) architectures will likely have multi-core nodes and may have more than one GPU. The nuances created by these variations to the primary design configuration are too varied to completely consider in this effort.

1) *Intranode Implementation:* GPU computing, as currently implemented, relies on massive data parallelism to realize speed-up in code run and compute times. Implementation is not without its challenges, as the programmer is restricted to problem domains which only allow for the data element to be treated in isolation or with minimal coupling to other data elements<sup>2</sup>. In particular, GPUs have emphasized high bandwidth local memory of limited size but at the cost of data transfers both across the PCI bus and from global GPU memory to fast local memory. Additionally, many of the tools and optimizations available to the CPU—such as optimized  $O(n \log n)$  variants of naïvely  $O(n^2)$  algorithms [8]—are not

<sup>1</sup>More precisely, we're mentioning different spatial scales. For multi-physics models requiring variable timesteps across the mesh, cell-based AMR can be used, but new challenges arise to match fluxed quantities across cell boundaries. Quirk [4] mentions specifically the issues of spurious reflections, but also provides an important discussion of grid efficiency and vorticity generation in adaptive meshes.

<sup>2</sup>This is not to say that problems of irregular and data-dependent parallelism have not been addressed. A review of some of the works at the GPU Tech Conference, <http://www.gputechconf.com/gtcnew/on-demand-gtc.php>, such as S0600, S0314, and S0042, show how graph representations and efficient implementations of scan/reduction operations can be utilized to see impressive performance. The hash-based algorithms developed by two of the current authors in Robey *et al.* [7] exploits a customized scan operation for speed-up. However, to get the best performance from the GPU, a high level of data independence is imperative.

available on the GPU<sup>3</sup>.

Nevertheless, GPU-based programming remains attractive to developers of scientific and numerical applications due to the extremely parallel nature and high memory bandwidth of the GPU; in particular, tasks such as numerically intensive linear-algebraic calculations can be executed in a significantly reduced period of time as compared to the same calculation performed on the CPU.

This speed-up is due to the large number of threads brought to bear by the GPU along with the essentially zero context switching time between the thread groups. However, given the evolving nature of an adaptive mesh, dynamical data structures are required which consequently cannot be solely handled by the GPU. More correctly, memory cannot be dynamically allocated on the GPU as of OpenCL specification 1.1 [1]. But this is easily solved; memory management on the CPU and numerical calculation on the GPU are combined to form a heterogeneous computing environment. It is here that the necessity of locality becomes apparent: there is a 20–40× factor increase in clock cycles due to memory latency when comparing read operations from global memory to read operations from local memory on the GPU [9, Ch. 3], with an even larger factor slowdown resulting from write operations back to the CPU [10, Ch. 3]. It is perhaps best to view the local memory on the GPU as a programmable cache for which the speed-up is highly dependent on the effectiveness of the programmer’s reuse of data while it is in the local memory.

2) *Internode Implementation:* For our single node computation, all calculations, including physics calculations, global reductions, neighbor calculations, and cell refinements, were successfully moved over to the GPU, with the CPU merely acting as a mechanism to reallocate memory. Namely, the state variables of the cells are resident on the GPU. As we move to MPI, however, our simple one node platform of one CPU core communicating with one GPU device needs to be expanded. The memory on the GPU must be retrieved by the CPU to communicate among nodes. Future enhancements to OpenCL will likely allow device buffers to be sent directly by MPI, but the complexity of fully implementing this technique will be challenging and require dynamic memory lists on the GPU.

The move to MPI sees an increase in the number of nodes, with each node still defined as a single CPU core controlling a single GPU device. We prefer a single core since our algorithm relies on the full use of the GPU for the numerical methods and additional CPU cores are a power drain that don’t contribute much to our performance. The GPU is so much faster than the CPU that it makes more sense to do all possible computation on the GPU. Still, most future CPU architectures will likely be multicore. While the additional CPU cores give us no advantage for a single node, since all the data is resident on the GPU and very little is done host-side on the CPU, data will need to be pulled off the GPU for MPI communication and writing out to files. The additional

CPU cores can then be used to speed-up these operations and provide for another level of data locality with some benefits of better cache use. To take advantage of this additional level of data hierarchy, we would assign a rank to each CPU core. Unfortunately, the current CUDA implementations give much worse performance with multiple CPUs driving one GPU. Thus, for this work, on each node, two MPI ranks were used to drive two GPUs on each node. Future CUDA implementations already have enhancements called Hyper-Q for multiple command queues on a GPU from multiple CPUs driving a single GPU. The effectiveness of this enhancement for multiple command queues cannot be evaluated until Kepler GPUs are available.

There are alternatives to this compute model. One would be to add an OpenMP or thread layer and have a single MPI rank for each node and use the CPU threading to gain access to the additional CPUs on the node. While technically viable, the essentially three different levels of parallel coding adds more complexity than we wanted to tackle. Another approach would be to run OpenCL kernels on the multicore level to access the additional CPUs. This option might be attractive in the future, but currently the GPU is so much faster than the CPU that it is hard to find work that would be reasonable to put on the multicore CPUs. The most attractive alternative at this time is to have one MPI rank per GPU and if there are additional CPU cores, just not use them in the computation. There is some waste in this compute model, but most large production systems will have a small number of CPU cores per node to conserve power. The last alternative would be to use CUDA on the GPUs instead of OpenCL. We have chosen OpenCL to allow a wider range of system architectures with some loss of additional CUDA functionality that could give higher performance in the near-term. Future hardware developments along the lines of the Intel MIC and AMD Fusion approaches may require a reassessment of our compute model, but it is clear that no single compute model will be able to run efficiently on every system architecture in the near future as hardware designs proliferate.

Our considerations in decomposing the mesh are focused on a viable scheme for a two-level partition. With the addition of the MPI internode layer we need to maximize locality in order to reduce data transfers across compute elements. This is the crux of CLAMR’s design. Moving to MPI challenges the versatility of our choice of partitioning. A standard method for partitioning in MPI, using a recursive bisection, then combined with a standard lexicographic data order aligned with cache on individual nodes, vastly complicates code complexity in load-balancing and MPI communication in general. Maintaining a data order to be used across all levels of the hierarchy, however, simplifies implementation while simultaneously giving the algorithm designer freedom to choose a data order which is most beneficial to the problem at hand. By using space-filling curves, global arrays can be partitioned across MPI compute elements by simply dividing the arrays by a fixed stride, ensuring balanced load balance; within the same data order, each MPI compute node can repeat this process of dividing the arrays by a fixed stride in order to best achieve load balance across the workgroups on the GPU.

<sup>3</sup>This is in reference to many of the tree-based algorithms which attain the  $O(n \log n)$  optimizations. Many tree-based algorithms are being ported to the GPU, but code complexity begs the question is the product worth the effort when a careful consideration of the problem decomposition might reveal underlying structure of a more independent nature.



### III. ARCHITECTURE & ALGORITHMS

The examination of heterogeneous platforms, discussed in Section II, led to several key architecture decisions for CLAMR. Our implementation is designed to efficiently create a dynamic memory space by using the CPU to manage memory transfers while the GPU performs the operations on the cells; memory transfers are reduced to maximize the time for which the control flow stays with the GPU. Calculating global reductions on the GPU is one technique used to accomplish this feat.

This section proceeds as follows. First an overview of the control flow of CLAMR is given. Second, the physical model and numerical method are analyzed, pointing out the issues involved in finite differencing across adaptive meshes. In the last three subsections, partitioning and locality, neighbor searching, and enhanced-precision sums, we look at the tools required to put cell-based AMR on the GPU. Specifically, data ordering is treated as a free parameter, allowing a “surface area to volume”-type optimization in array distribution between processes—and thus a minimization in the amount of data requiring transfer. As a consequence, neighbor searching becomes nontrivial and global calculations vary, hence the necessity of clever neighbor searching methods and consistency checking enhanced-precision sums.

CLAMR is available through an open-source license from [github.com/losalamos/CLAMR](https://github.com/losalamos/CLAMR).

#### A. CLAMR: Control Flow

The initialization of CLAMR begins by creating global objects. The mesh is built with each cell’s initial state variables set to the problem specification. This is followed by a preliminary global space partitioning accompanied by a calculation of cell neighbors. Next, the compute context is established, which includes a command queue containing the commands to be sent to the compute device. In particular, all kernel objects are declared. (For a review of OpenCL terminology, see the OpenCL Specification [1].) As part of establishing the compute context, memory is allocated on the GPU. Then the state variables for each cell are written to the GPU’s global memory space, and all kernel arguments are set.

When setup completes, control flow is transferred to the GPU as the size of the timestep for the numerical scheme is computed. This is calculated based on the wave speed which requires a global reduction as the maximum wave speed needs to be established. Then, in preparation for execution of the workgroups, the local tile (workgroup memory space variables) is set. In addition, the conditions on the outer boundary of the mesh are created as needed.

At this point, the state variables are updated as determined by the governing equations and numerical scheme. Section III-B provides the details. We note here that calculations are done in double precision on both GPU and CPU. Following the update, the gradients are used, depending on the magnitude and sign change across cells, to refine or coarsen the cells of the mesh. A device-global reduction is done to indicate the new number of cells, both globally and on each tile.

With the new number of cells in the mesh determined, memory on the GPU needs to be reallocated. Control flow is returned back to the CPU; the new variable arrays are first resized before the rezone call and the old arrays are resized afterwards in a technique reminiscent of the double-buffering commonly done in graphics applications. Control flow is once again returned to the GPU as a rezoning of the cells is done; cells are refined, coarsened, or unchanged as necessary. The space-filling stencil is applied (which is highly parallel, namely a global space-filling curve call is no longer needed), and the new cell neighbors are set.

This procedure is iterated for the desired length of the simulation.

#### B. The Physics Model & The Numerical Method

Currently, CLAMR implements the shallow water wave equations because of their relative simplicity as well as the high degree of symmetry present in our particular initial condition: a cylindrical shock impacting the center of the mesh (*i.e.*, a circular dam break). In their conservative form, the equations are:

$$\begin{aligned} \frac{\partial h}{\partial t} + \frac{\partial(hu)}{\partial x} + \frac{\partial(hv)}{\partial y} &= 0 \\ \frac{\partial(hu)}{\partial t} + \frac{\partial}{\partial x} \left( hu^2 + \frac{1}{2}gh^2 \right) + \frac{\partial}{\partial y} (huv) &= 0 \\ \frac{\partial(hv)}{\partial t} + \frac{\partial}{\partial x} (hvu) + \frac{\partial}{\partial y} \left( hv^2 + \frac{1}{2}gh^2 \right) &= 0, \end{aligned}$$

respectively conservation of mass, conservation of  $x$ -momentum, and conservation of  $y$ -momentum. Here  $h$  is the height of a column of water,  $g$  is the acceleration due to gravity, and  $u$  and  $v$  are the wave velocities in the  $x$  and  $y$  directions, respectively. More precisely, they are the velocities of the water molecules. The speed of propagation, the phase velocity, is  $\sqrt{gh}$ . For a nice discussion which mentions this, see Tao [11]. Note that mass equals height (times a constant) because water is incompressible. In particular, the incompressibility causes all pressures to only vary the height of the water while the width and length of a differential column remain constant, and the density of water remains constant. Also note that the pressure term is  $gh^2/2$ . For a more rigorous presentation, see the sections in Landau & Lifshitz [12] on long gravity waves and shallow-water theory. If the reader is particularly interested in the shallow water equations on the GPU, we suggest Castro *et al.* [13] and Brodtkorb *et al.* [14].

For the discretization, we use a total variation diminishing (TVD) finite difference scheme based on a two-step Lax-Wendroff method [15] in conjunction with a minmod symmetric flux limiter to provide an upwind weighted artificial viscosity term. The Lax-Wendroff method is second-order accurate in space and time, hence it is a suitable choice for smooth regions. Around steep shocks, oscillations produced by the second-order method necessitate damping by switching to a first-order upwind method, which is accomplished by the

flux limiter. For a solid foundation on the numerical methods mentioned, as well as modeling the shallow water equations in general, see LeVeque's book on finite difference methods [16] and George's Master's thesis [17]. For more information about the complete numerical method as implemented, see Davis [18], Sweby [19], and Yee [20].

While the references above provide the background for the numerical method used in CLAMR, the adaptive mesh complicates the equations. The half-timestep equations are:

$$U_{i+1/2, j}^{n+1/2} = \frac{r_i U_{i+1, j}^n + r_{i+1} U_{i, j}^n}{r_{i+1} + r_i} - \Delta t \left( \frac{F_{i+1, j}^n A_{i+1} a_{i+1} - F_{i, j}^n A_i a_i}{V_{i+1} v_{i+1} + V_i v_i} \right)$$

$$U_{i, j+1/2}^{n+1/2} = \frac{r_j U_{i, j+1}^n + r_{j+1} U_{i, j}^n}{r_{j+1} + r_j} - \Delta t \left( \frac{G_{i, j+1}^n A_{j+1} a_{j+1} - G_{i, j}^n A_j a_j}{V_{j+1} v_{j+1} + V_j v_j} \right).$$

Following a standard notation, subscripts  $i$  and  $j$  are spatial indices, while the superscript  $n$  is a time index. Here  $U$  represents a general state variable, which for the shallow water equations are the mass, the  $x$ -momentum, and the  $y$ -momentum.  $F$  and  $G$  are the  $x$  and  $y$  flux terms for the state variable  $U$ , respectively.

The adaptation from the regular grid equations,

$$U_{i+1/2, j}^{n+1/2} = \frac{U_{i+1, j}^n + U_{i, j}^n}{2} - \frac{\Delta t}{2\Delta x} (F_{i+1, j}^n - F_{i, j}^n)$$

$$U_{i, j+1/2}^{n+1/2} = \frac{U_{i, j+1}^n + U_{i, j}^n}{2} - \frac{\Delta t}{2\Delta y} (G_{i, j+1}^n - G_{i, j}^n),$$

arises from the inclusion of spatial scaling variables which are necessary to account for the adaptive mesh. In the adaptive mesh half-timestep equations,  $r$ ,  $A$ , and  $V$  are radius, area, and volume, respectively, along with the scaling variables  $a$  and  $v$  for the area and volume, respectively. The first expression with the  $r$  terms is a linear interpolation to find the state variable at the face between cells (the large  $\times$  between cells 5 and 7 in Figure 1). In the regular grid equations, this interpolation is merely the average of the two neighboring state variables. The second expression which includes the timestep is the flux term; it weights the fluxes  $F$  and  $G$  by the area of the cell interface multiplied by a scale factor, which to allow for arbitrary dimension is taken in the form

$$a_i = \min \left( 1, \frac{a_{i+1}}{a_i} \right), \quad a_{i+1} = \min \left( 1, \frac{a_i}{a_{i+1}} \right).$$

This total flux, expressed in the numerator of the second term of the adaptive mesh half-timestep equations, is then divided over the volume of the staggered compute cell (the colored region in Figure 1), which appears in the denominator. In order to allow for arbitrary dimension, the volume contributions by neighboring cells are multiplied by the scale factors

$$v_i = \min \left( \frac{1}{2}, \frac{v_{i+1}}{v_i} \right), \quad v_{i+1} = \min \left( \frac{1}{2}, \frac{v_i}{v_{i+1}} \right).$$

In the regular grid equations, no scaling is necessary, and consequently  $A/V$  is simply  $1/\Delta x$  or  $1/\Delta y$ .

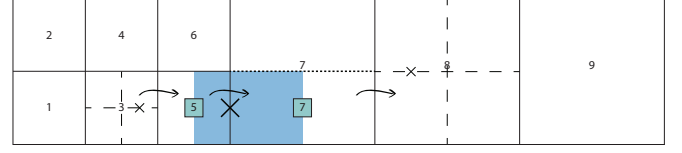


Fig. 1. The finite difference stencil shows the half-timestep calculation, as well as the comparison of gradients for the flux correction to remove oscillations (applied at the full-timestep). The value of a state variable is linearly interpolated at the face (linear interpolations are shown with an  $\times$ ), and the fluxes for the staggered compute cell are computed from the stored values (green boxes). Notice that the value of cell 7 is taken along a half cell with the same characteristic width as cell 5. That is, the neighbor of greater refinement sets the characteristic scale. When flux limiting is applied, a five-point stencil is used to find sequential gradients (shown by the arrows).

Efficient use of the GPU requires that the equations be consolidated into as few as possible so that all cells/threads execute the same code block. The same equation expressed in multiple conditional blocks will suffer performance degradation due to all threads in a workgroup having to execute each block. That is, breaking an equation into two similar blocks controlled by a conditional statement—such as if the cells differ in refinement level—will effectively double the runtime since all threads in the workgroup execute both blocks (although only results from the appropriate operations are preserved). This characteristic of GPU performance has been termed “lock-step”. Hence consolidating the equations into one with additional factors to scale or turn on/off terms is crucial for GPU-based code performance. These extra factors are calculated at the start of the kernel in short conditional blocks, thereby minimizing the consequent performance reduction. Reference [21] discusses issues of divergence and lock-step on the GPU.

The full-timestep calculation is:

$$U_{i, j}^{n+1} = U_{i, j}^n - \Delta t \left( \frac{\overline{F}_{i+1/2, j}^{n+1/2} - \overline{F}_{i-1/2, j}^{n+1/2}}{\Delta x} + \frac{\overline{G}_{i, j+1/2}^{n+1/2} - \overline{G}_{i, j-1/2}^{n+1/2}}{\Delta y} \right).$$

The important point to note here is that the calculation from the perspective of a coarse cell with refined neighbors requires averaged fluxes, denoted by the  $F$  and  $G$  terms with bar overhead. This averaging, however, is complicated by the flux terms for one state variable not being of the same form as another. For example, the mass flux term in the  $x$  direction is  $hu$ , while the corresponding term for the  $x$ -momentum is  $hu^2 + \frac{1}{2}gh^2$ . Looking again at Figure 1, we see that the flux across cell 7's left face is the sum of the fluxes across its interfaces between cell 5 and cell 6. Numerically the most important consideration is conserving the state variables, and this equates to computing the fluxes identically from the perspectives of cells 5, 6, and 7. The major ramification is the necessity of averaging the fluxes at the full-timestep, as opposed to averaging at the half-timestep.

Finally, to correct for oscillations near shocks, a minmod symmetric flux limiter is used to impose the TVD property.

This requires a five-point stencil in order to take state information from the neighboring cells' neighbors, thereby allowing the ratios of the gradients to be examined. This of course is in contrast to the above equations for the Lax-Wendroff method, which are compact in the sense that they are only using a cell's nearest neighbors' state information. Then, once the gradients are compared, and if there is a sign change, the flux limiter term is applied. The corrections are:

$$U_{i,j}^{n+1} \pm \frac{\nu(1-\nu)}{2} [1 - \phi(r^+, r^-)] \Delta U^n$$

$$\phi(r^+, r^-) = \max(0, \min(1, r^+, r^-))$$

where  $\phi$  is the flux limiter, and  $\Delta U$  is an upwind difference in the state variable. The Courant number,  $\nu$ , is the timestep divided by the grid spacing and multiplied by an eigenvalue of the system of equations corresponding to the state variable. The ratios  $r^+$  and  $r^-$  are dimensionless values quantifying the change in the gradient across the five-point stencil. Referring to Figure 1, where the flux correction is being computed at the interface between cells 5 and 7,  $r^+$  and  $r^-$  are determined by taking the inner product of sequential finite differences and dividing by the finite difference at the interface. For example,  $r^+$  takes the inner product of the gradient across the interface between cells 7 and 8, as shown with the arrow (and which may or may not require an interpolation as shown by the small  $\times$ ), with the gradient across the interface between cells 5 and 7. This is then divided by the gradient across the interface between cells 5 and 7, squared, thereby effectively capturing the sign change in the upwind direction. The value of  $r^-$  is computed in the same manner, only replacing the gradient across the interface between cells 7 and 8 with that between cells 3 and 5.

### C. Partitioning & Locality

The numerical method presented above was formulated in such a way as to take full advantage of the maximum throughput of the GPU. The physical calculation is compressed into a few equations with minimal conditional branching, allowing all processing elements in a GPU workgroup to efficiently cooperate. However, this only ensures that a work-item processing data for a single cell can maximize its concurrency with other work-items in the workgroup. An efficient scheme for partitioning cells, and therefore computations, homogeneously across compute cores is still absolutely imperative.

Computer data representation of a multidimensional array is, of course, actually linear, with data offsets calculated as a function of the row and column of the index; thus, taking an arbitrary  $k$ -dimensional mesh and mapping it into a one-dimensional array provides a direct, easily applied method for apportioning data elements across the GPU's cores. Elements of this 1-D array can be taken sequentially in blocks sized to match the workgroup memory. Even with this efficacious utilization of available memory, limited workgroup memory and relatively expensive calls to global memory provide an incentive for keeping as much data as possible as local as possible. There is a strong motivation for ordering the array such that locality is largely preserved, as this eliminates cache

and page misses (access operations beyond the scope of the current cache or page). Accordingly, a key goal of the data structure used for CLAMR is to decompose the two-dimensional grid of state variables into a linear array while minimizing the number of out-of-workgroup neighbor accesses that must be made.

Judicious ordering of mesh elements in the one-dimensional array can mitigate this problem. Consider two different contiguous load divisions of a two-dimensional grid of work-items between several workgroups, each capable of processing  $N$  work-items (Figure 2). In the naïve contiguous case, ordering the work-items linearly by column and row requires  $2N + 2$  off-tile accesses per workgroup (taking all  $N$  elements to be in a single row as the characteristic case). The optimal case for a perfect square tile only necessitates  $4\sqrt{N}$  off-tile accesses. It is easy to see that this problem is basically a "surface area to volume" ratio minimization, requiring that the mesh preserve local structure rather than the straightforward ordering. Also, it has been found that the addition of refinement to the problem aggravates the non-locality problem by progressively fragmenting the mesh representation for work-items in a workgroup if cautious measures are not taken.

To best achieve surface area to volume minimization, the use of appropriate space-filling curves for filling a  $k$ -dimensional space is necessary. Peano [22] explored the theory behind mapping higher-dimensional spaces to a one-dimensional representation (with the same cardinality) by using a self-recursive stencil<sup>4</sup>. For the purposes of the current exposition, a comparison of the Hilbert [25] and Z-order [26] space-filling curves suggests that the Hilbert curve maximizes spatial locality (Section IV).

A true space-filling curve recursively enumerates uniform divisions of a space at higher and higher levels of refinement; it is straightforward to see how this can be carried over to a mesh with multiple levels present simultaneously. The global Hilbert curve requires a recalculation of the full index every time a refinement occurs anywhere in the mesh; the practical application of this technique on a multilevel grid can lead to some unusual dead ends and consequent backtracking in the index ordering (Figure 3), as well as providing a code bottleneck<sup>5</sup>. Consequently, a local stencil preserving the entry and exit neighbors before refinement has been implemented, guaranteeing locality while avoiding fragmentation and backtracking/multiple indexing. The relative orientation and direction of the required coarse curve through the cell is calculated, and a unique local ordering is selected via a rule-based scheme matching the Hilbert curve of the refined mesh which preserves itself through arbitrary mesh refinements and coarsenings. A restructuring of the calculation of the Hilbert curve based on geometric location rather than entry/exit attributes can incorporate this local stencil concept.

Load-balancing, one of the important considerations in distributed-memory message-passing systems, is cleanly

<sup>4</sup>see Figure 2 of Haverkort *et al.* [23] for a visual survey of several space-filling curves and Jin *et al.* [24] for a discussion of space-filling curves as a means for computational reordering

<sup>5</sup>Also, note that repeated indices do not mean repeated cells, only that there are multiple referents to that cell in the neighbor array.

solved by this implementation on the GPU. It can spawn a number of work-items not restricted by the number of physical processing elements on the device. Coupled with a strategy for ordering the indices which minimizes the surface area of the tile, this scheme distributes the processing and memory access load equitably and nearly optimally (For further internode discussion, see Section II-B2).

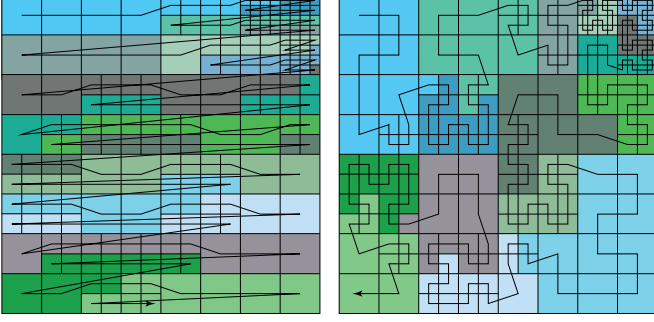


Fig. 2. Here are two alternate load divisions of a 2-D array between several workgroups (indicated by color). On the left is a naïve division which begins to fragment with refinement, requiring  $\sim 2N + 2 \rightarrow O(N)$  off-tile accesses for an  $N$ -element tile. On the right is a more optimal division done by the Hilbert curve, which nearly minimizes the surface area to volume ratio, consequently requiring only  $\sim 4\sqrt{N} \rightarrow O(\sqrt{N})$  off-tile accesses for an  $N$ -element tile.

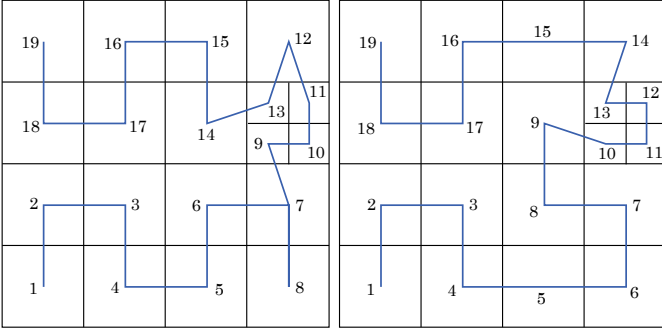


Fig. 3. The Hilbert space-filling curve for a coarse mesh is shown locally refining with a global curve calculation on the left and locally refining with a local curve calculation on the right.

#### D. Neighbor Searching

The price to be paid for an optimal spatial partitioning of the mesh is a nontrivial order of the data elements. This aggravates the neighbor calculation and makes the use of the basic stencil of the numerical method nontrivial. For instance, for the Lax-Wendroff method, only a cell's most immediate neighbors' state information is needed; the TVD correction algorithm requires access to second-nearest neighbors.

From the hardware perspective of the GPU, we know that workgroups cannot communicate with each other and that the access times to global memory are much slower than the access times to local memory. In order to take advantage of the spatial locality achieved by the partitioning of the last section, efficient search and retrieval of neighbor data is essential. All necessary state information should be preemptively stored in local variables at the start of the physics calculation kernel.

One approach is to pass neighbor information onto the GPU in arrays of neighbors. Due to the optimization in

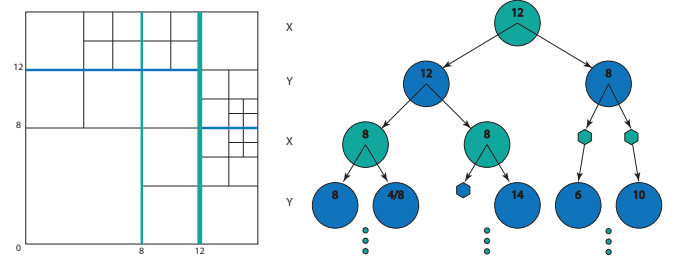


Fig. 5. The mesh is recursively bisected using a  $k$ -D tree.

locality achieved by the use of the Hilbert curve, a significant proportion of the work-items would be able to make local GPU memory accesses for neighboring cell information, although a few would still need global memory accesses. These latter cases induce conditional branching in the main computation kernel. Additionally, only one neighbor in each direction is stored in this scheme, requiring two neighbor accesses to access all data across refinement levels (Figure 4). (*i.e.*, in order to compute the flux across a face for which the bordering cells are at a level of greater refinement, two neighbor accesses are required for a single face. However, the current cell doesn't have direct access to one of those neighbors, and so it must first query the neighbor it does know to get the index of the unknown neighbor in the Hilbert-curve-ordered array.)

Constructing the arrays of neighbor indices could proceed in several ways. A very naïve algorithm would simply take the current cell and search through every other cell of the mesh to check whether it's a neighbor, *i.e.*, an  $O(n^2)$  algorithm. A more clever approach uses a  $k$ -D tree as seen in Figure 5, which recursively bisects the mesh using a weight function (which in this case is the number of cells) while alternating axes. Construction of the arrays of neighbors then requires  $O(n \log n)$  time. A skip-list is a further optimization utilizing a hierarchy of linked-lists (effectively a hierarchy of trees with various levels of sparsity), thereby introducing speed-up as a result of random accesses; unfortunately, the worst case can still give  $O(n \log n)$  time complexity. Nevertheless, there is active research in construction of  $k$ -D trees. (Reference [27], Section 4.1, discusses building a  $k$ -D tree on the GPU.)

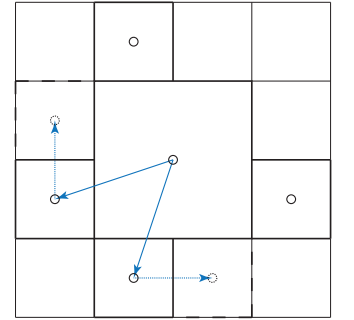


Fig. 4. The neighbor accessing scheme illustrates the need to access the neighbor of a neighbor given certain levels of relative refinement.

In practice, it was found that construction of a  $k$ -D tree at every timestep introduced a significant performance bottleneck. The construction of a  $k$ -D tree on the initial timestep, followed by regular updating of the neighbor indices by tracking refinements and calculating relative offsets, introduced significant complexity into the code as well as implicitly assumed a Z-order space-filling curve. A better (and faster)



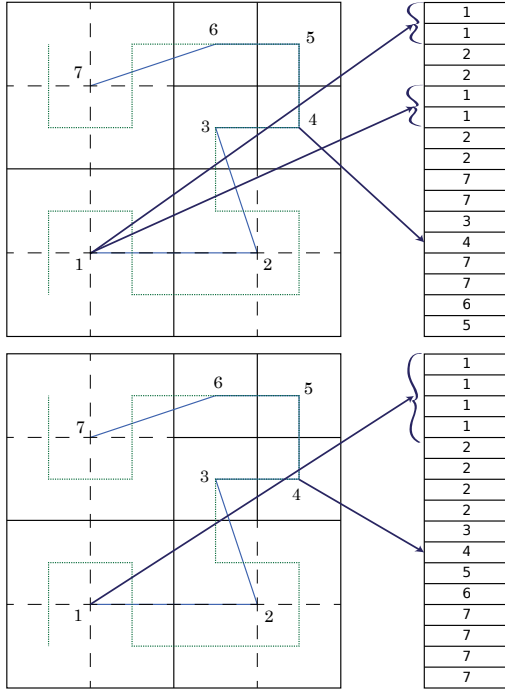


Fig. 6. The first hash mapping follows the standard indexing used for arrays, spanning the  $x$ -axis first and then incrementing the  $y$  coordinate after one stride has finished. This is done on a superimposed grid of the finest level of refinement. The second hash mapping macroscopically matches the indexing induced by the Hilbert space-filling curve for a superimposed grid of the finest level of refinement.

solution is to use an analytic function to map cells into a hash table, requiring only  $O(n)$  time for a linear-time look-up, with the added advantage that the hash table construction is done on the GPU. The design target was to find a perfect (collision-free) hash for simplicity and performance reasons. Exploiting the underlying structure of the AMR grid and that no point in the mesh could exist at two levels at one time, the minimal memory for a perfect hash table would be the size of the grid at the finest allowed mesh refinement, thus allowing the mapping of all possible cells with no extra space incorporated. As the maximum number of levels of refinement is increased, however, the hash table size can grow quickly, as discussed in Section IV-D. Figure 6 depicts two alternate hash schemes, with Figure 7 providing a reference for the indexing. For more in-depth work on parallel hashing with the GPU, including collision handling algorithms, see Alcantara *et al.* [28].

1) *Hash Implementation and Key Functions:* The hash indexing currently implemented uses a standard Cartesian coordinate defined key. A cell, regardless of its refinement, calculates its position as if it were at the finest level of refinement, and then it can make a constant time access to the hash table. Its position in the hash table is indexed as  $w \cdot y + x$ , where  $w$  is the width of the mesh, and  $y$  and  $x$  are its Cartesian coordinates, all taken at the finest level of refinement.

In the CSAMR code [29], a clever oct-based data structure is presented for an adaptive mesh. A hash function was used to access parts of the mesh, but each level of refinement was also stored in the hash table, being more prodigal with memory

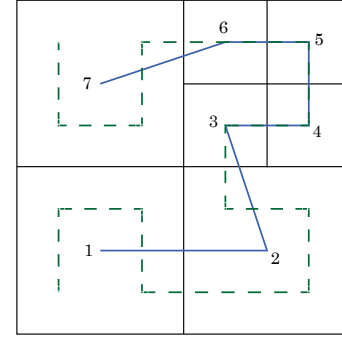


Fig. 7. This is the reference grid for the hashing done in Figure 6. The blue curve fills the actual mesh, while the dotted green curve shows how the mesh would be filled if it were entirely refined at the finest level.

than the method implemented here. For the GPU, CLAMR needs to minimize the amount of memory required; thus, the hash table is structured at the greatest level of refinement and coarser cells merely occupy multiple elements in the array. In CLAMR, there is actually an optimization which guarantees a cell never writes to more than 7 hash buckets, a byproduct of the rule that neighboring cells can only differ in scale by one level of refinement (depicted in Figure 6).

On the other hand, given the standard indexing, there is significant divergence from the Hilbert curve in a cell's location in the arrays. As an interesting course for exploration, one could use the hash mapping to tandemly structure the hash table in order to match the Hilbert-curve partition of the finest level of mesh cells. Gains are more likely to be seen here than in the use of MPI, as this will reduce transfers of the pieces of the global hash table.

#### 2) Algorithm Block for Neighbor Calculation Using MPI:

The biggest change associated with adding the MPI layer is in the neighbor calculation routine. The construction of the hash table necessitates focus on the local region of the mesh in order to achieve satisfactory memory scalability. The algorithm becomes more complicated as shown in Algorithm 1.

This is the most complicated and error-prone portion of the MPI/OpenCL code. However, once this is done, the rest of the code simply uses the ghost cell update calls with little additional complexity over the serial code.

#### E. Enhanced Precision Sums

Modifying the data order has the undesirable side effect of changing the calculation of the total mass in a problem by a small amount ( $\sim 10^2$  the machine epsilon). This is due simply to the finite-precision floating-point arithmetic of adding up the elements of a large array for which addition is not truly associative; though well within the acceptable error of the simulation, it complicates debugging and verification of the programming. The problem is exacerbated as the problem size and data range increase. Thus, to properly allow arbitrary data reordering for performance and locality, this source of error should be dealt with by the use of an enhanced precision sum.

Enhanced precision sums can be used to reduce or even eliminate global sum errors. The basic concept is to use a second variable to carry the truncation part of the sum thereby



**Algorithm 1** Neighbor Calculation Using MPI

- 1) Determine min/max  $i$  and  $j$  of local mesh region.
- 2) Add 2-cell extra buffer of coarse “ghost cells” at the edge of the mesh to allocate hash table.
- 3) Calculate local hash table.
- 4) Calculate local neighbors.
- 5) Determine all cells on the boundary for layer 1 of the local tile by searching for unsatisfied neighbors.
- 6) Look inward from boundary cells for index data for layer 2 needed by other processors and gather to all processors.
- 7) Each processor fills in off-tile extra buffer with gathered values.
- 8) Find layer 1 of off-processor data needed from new neighbors.
- 9) Find layer 2 of off-processor data
- 10) Fill-in ghost region in arrays
- 11) Recalculate neighbors for ghost regions with new hash.
- 12) Set up communication pattern for ghost regions.

increasing the digits of accuracy. Prior studies (Robey *et al.* [30]) showed that the Kahan enhanced precision sum operation reduces the error in calculation as well for distributed MPI-based methods.

A demonstration problem was developed to stress the enhanced precision sums with a million cells, six orders of magnitude range in the initial conditions, and 3 levels of refinement. Every time the mesh was refined, a global sum was done for the Z-order and Hilbert data orderings. The maximum difference between the global sums for the two data orders in 100 iterations was measured at 448 times the machine epsilon. This is high enough to mask small errors in programming such as using old boundary conditions. Using the Kahan sum, the maximum difference in the sums for the two data orders was reduced to zero, vastly improving the detection of programming errors.

The cost of this technique is so small that it should be standard practice, particularly when data orderings are changing as is commonly the case in parallel calculations. The cost of performing the sum on the CPU is about 66% over that of the standard sum. Since the sum is less than one percent of the run-time cost, the addition to the run-time is negligible.

#### IV. PERFORMANCE AND SCALABILITY

With the use of data ordering as a free parameter, the GPU-enabled code achieved an overall speed-up factor of one order of magnitude over the CPU-only code for our single node runs. Tests of CLAMR were performed on the Moonlight cluster at Los Alamos National Laboratory. A node on the cluster is comprised of two eight-core Intel Xeon E5-2670 CPUs rated at 2.6 GHz, each core with 0.5 MB of secondary cache and each chip with a 2MB tertiary cache between its eight cores. The two Xeon sockets share 32 GB of RAM on each node. Each compute node has two GPGPU NVIDIA Tesla M2090 cards connected to PCIe-2.0 x16 slots. Practical maximum bandwidth to these GPGPU cards is between 6.0 GB/s to 6.6

GB/s<sup>6</sup>. For the test runs, however, a single node is defined as using a single core of one of the Xeon chips, and using a single Tesla M2090. CLAMR was compiled and linked using GCC 4.4.6 as the C compiler and the NVIDIA OpenCL 1.1 SDK from the Cuda toolkit v. 4.1. Tables I and II provide a detailed breakdown of the dominant routines in CLAMR on a  $256 \times 256$  coarse mesh and a  $512 \times 512$  coarse mesh, respectively, with 2 and 4 levels of refinement. Explicit CPU use is compared to the hybrid implementation of the CPU combined with the GPU. Note that the mesh does not refine every cycle; for the  $256 \times 256$  mesh at 2 levels it refined 28% of the cycles and for 4 levels it was 65%. For  $512 \times 512$  coarse mesh with 2 levels refinement occurred 45% of the iterations and for 4 levels it occurred 87% of the iterations.

The significance of the results stems not from the absolute performance gain but from the realization that getting on the GPU performance curve is vital to future performance increases. For example, earlier runs were done on a cluster where each node was comprised of an AMD Opteron 6168 processor for the CPU and an NVIDIA Tesla C2050 for the GPU. Those runs showed an overall  $35\times$  speed-up. Likewise, as the NVIDIA Kepler cards are released, the speed-up factors should increase again. Hence an order of magnitude speed-up is well worth the effort, especially as the GPU performance curve promises to drive this factor higher. Nevertheless, the numbers show the authors that there are still further optimizations desired in the finite difference kernel, as well in the hash setup of the neighbor calculation.

##### A. Partitioning Data

Demonstrating quantitatively the effectiveness of various partitioning via space-filling curves, Figure 8 shows that the Hilbert curve achieves a nearly optimal surface area to volume ratio minimization. The number of out-of-cache accesses for a nearest-neighbor algorithm is governed by the number of bordering tiles, and is thus minimized for a contiguous array by reducing the perimeter-to-area ratio, which has a minimum of  $4\sqrt{N}$  (for square  $N$ ). Hungershofer and Wierum [31] propose the normalization of the perimeter-to-area ratio for different partitions by dividing by the minimum perimeter-to-area for the square or the cube in three dimensions. Normalizing the measure makes comparison of different partition sizes easier since a value of 1 is ideal. This partition quality coefficient, or C value, is defined in two-dimensions as

$$C^{curve} = \frac{S^{curve}(N)}{4\sqrt{N}}$$

The measure of the Hilbert curve is in stark contrast to the measure of the Z-order curve, for which increasing the maximum level of supported mesh refinement actually leads to an *increase* in the ratio (more fragmentation). While the original order partitioning slowly minimizes the ratio further with increasing refinement, it is clear that the Hilbert curve achieves the best spatial partitioning.

The effect of the local stencil on partitioning is also seen in Figure 8, in which a global Hilbert curve filling the space

<sup>6</sup>Specifications courtesy of <http://hpc.lanl.gov/tlcc2home>.

TABLE I  
PERFORMANCE ON A SINGLE NODE – 1000 ITERATIONS ON A  $256^2$  COARSE MESH

<i>Function</i>	2 Levels of Refinement			4 Levels of Refinement		
	CPU (s)	GPU (s)	Speed-up	CPU (s)	GPU (s)	Speed-up
Total	31.84	2.91	10.9	57.94	5.32	10.9
Time Step	1.57	0.21	7.4	2.04	0.27	7.5
Finite Difference	27.57	1.87	14.7	37.65	2.39	15.7
Refine Potential	1.47	0.58	2.5	2.71	0.85	3.2
Rezone All	0.37	0.10	3.8	1.10	0.16	6.8
Neighbor Calc	0.79	0.15	5.3	14.34	1.65	8.7

TABLE II  
PERFORMANCE ON A SINGLE NODE – 1000 ITERATIONS ON A  $512^2$  COARSE MESH

<i>Function</i>	2 Levels of Refinement			4 Levels of Refinement		
	CPU (s)	GPU (s)	Speed-up	CPU (s)	GPU (s)	Speed-up
Total	124.00	9.76	12.7	223.60	18.11	12.3
Time Step	5.98	0.38	15.6	6.93	0.42	16.3
Finite Difference	101.39	6.92	14.6	121.29	7.69	15.8
Refine Potential	6.06	1.27	4.8	10.24	1.79	5.7
Rezone All	2.59	0.41	6.3	5.78	0.50	11.5
Neighbor Calc	7.70	0.76	10.1	78.97	7.69	10.3

on every iteration is compared to an initial Hilbert curve followed by two separate local stencils for refinement—a fixed Z-order stencil and the Hilbert stencil. The partition measure demonstrates a small but measurable positive change with the introduction of the local stencil. The use of the local stencil for the Z-order actually counters the fragmentation associated with the global Z-order. Given the ease of implementation for a local stencil, and the inherently parallel nature of its application, these results provide good evidence for their use.

### B. Single-Node Speed-up of Partitioning

Of further consequence, the partition measure acts as an indicator of the probable effect a partitioning method has on run-time. Our analysis shows increasing correlation as the maximum level of refinement is increased. On the other hand, our partition measure is first-order in the sense that we have only considered spatial locality, *i.e.*, the amount of data which is to be transferred. A secondary measure might account for the effects of hardware, such as exploiting quad-loads from cache on the GPU. Consequently, the original order still has impressive performance on some hardware when comparing to the best space-filling curve, the Hilbert curve with the local stencil, since a single off-tile access has the tendency to grab a neighboring cell which will be needed by another cell in the workgroup. Shown below in Figure 9 is a collection of normalized runs on the aforementioned cluster comprised of single nodes with the Intel CPU and NVIDIA Tesla C2090

GPU as a function of partition measure. As the graph shows, partition measure based on spatial metrics is a solid indicator of runtime, except in the original order case in which cache access is not accounted for.

Nevertheless, the performance of the Hilbert curve with local stencil matches the original order and the uniform interface of the space-filling curve approach makes coding much simpler, especially in the context of MPI, considered next.

### C. MPI Performance: Strong Scaling, Weak Scaling, and Partition Measure

We tested the effect of our Hilbert curve scheme on both strong scaling (varying only the number of processing elements for a fixed problem size) and weak scaling (varying the problem size only for a fixed number of processing elements), comparing it to a scheme which uses the original data order both intranode and internode. All data presented in this section were obtained from runs on the aforementioned Moonlight cluster. The strong scaling tests were run with a fixed coarse mesh size of  $1024 \times 1024$  at both 2 and 4 levels of refinement for 500 timesteps. Figure 10 only shows the data for the Hilbert curve method, as the original order produced nearly identical data and hence for clarity is excluded. Even at a few nodes the MPI-only scheme quickly outperforms the single node CPU, as expected. The hybrid MPI/GPU scheme has more overhead of pulling data off the GPU for MPI communication. The MPI-

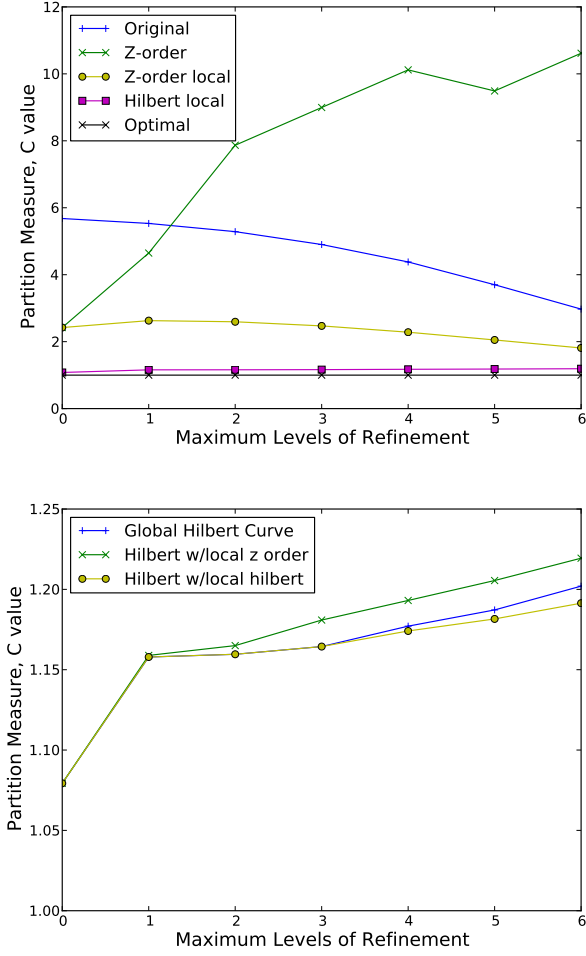


Fig. 8. The data plotted on top indicate that the Hilbert curve achieves a near optimal minimization of the surface area to volume ratio. Hence, fewer off-tile neighbor accesses are made. The plot on the bottom shows that the use of a local stencil has a measurable positive effect on the surface area to volume ratio.

only catches up to the MPI/GPU at about 64 nodes, but at that point there is only about 16,000 cells per processor negating the speed-up effects of the GPU.

The weak scaling tests were run with a base coarse mesh of  $512 \times 512$ , so that the problem size is  $512 \times 512 \times N$  where  $N$  is the number of nodes. Again the runs were done for 500 iterations at both 2 and 4 levels of refinement; the case for 2 levels of refinement is shown in Figure 11. We see in the top plot that, for both the Hilbert scheme and the original order scheme, the scaling efficiency is better for MPI only versus the MPI/GPU hybrid. This holds true for 4 levels of refinement as well. We also see that the scaling efficiency curves for both MPI and the MPI/GPU hybrid exhibit similarity; of particular interest is that the similarity is mirrored in the bottom plot of Figure 11 which shows the MPI partition measure as a function of the number of nodes. Here the partition measure is the surface area to volume ratio of a node's region of the global mesh. Figure 11 also shows that our Hilbert scheme obtains better scaling efficiency than the original order at 2 levels of refinement.

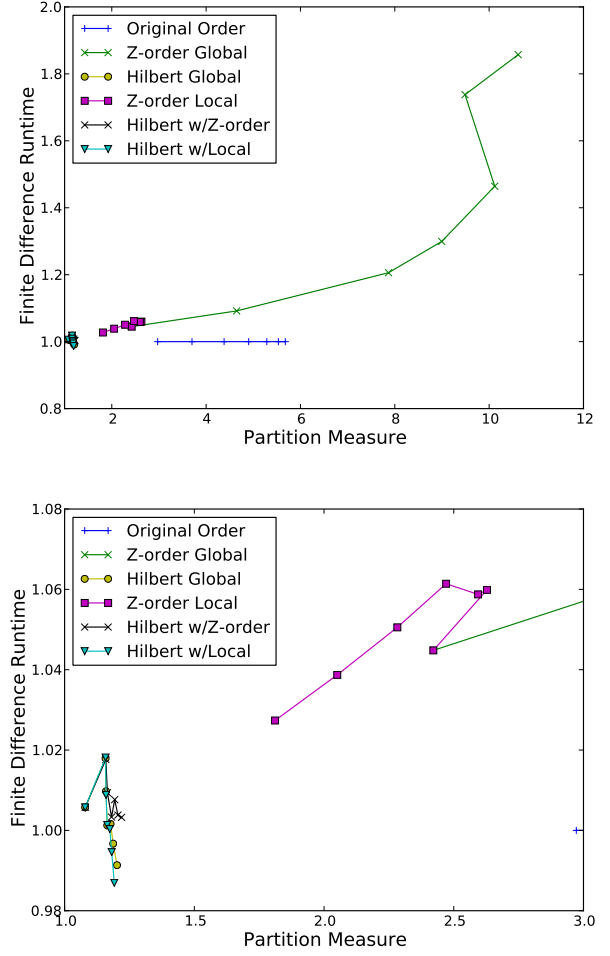


Fig. 9. Here is a collection of normalized runs as a function of partition measure. We see that the partition measure, based on spatial metrics, is a solid indicator of runtime except for the original order since cache accessing effects are not accounted for in the metrics.

Our conclusion is that the focus on locality is the right choice, given the similarity of the scaling efficiency curves to the MPI partition measure curves. Ultimately the partition measure needs to be expanded beyond solely “spatial” metrics and include the effect of cache behavior.

#### D. Neighbor Calculation Data

On a related note, the hash table implementation for neighbor calculations shows an impressive speed-up over the  $k$ -D tree, as seen in Figure 12. The CPU hash achieves an order of magnitude speed-up over the  $k$ -D tree, while the GPU hash achieves a speed-up of three orders of magnitude. By increasing the memory complexity of the neighbor algorithm, a large reduction from  $O(n \log n)$  to  $O(n)$  is achieved in the time complexity. There are also two desirable aspects regarding the additional memory space: it is used only temporarily, and each bucket in the hash table only stores a single integer for perfect hashing. In particular, production level codes can have hundreds of state variables per cell, each stored as a float or double, so the relative space taken in global memory for the hash table is not unreasonable.

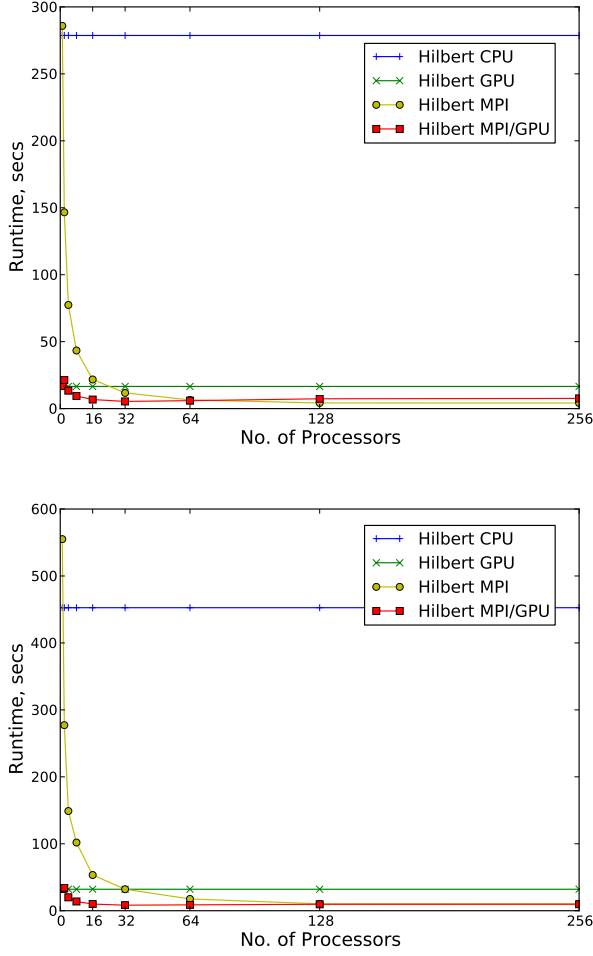


Fig. 10. The top graph shows strong scaling for a problem run on a  $1024 \times 1024$  coarse mesh with a maximum of 2 levels of refinement for 500 iterations, and the bottom graph shows the same problem but with a maximum of 4 levels of refinement.

Figure 12 does show that there is still an issue with the growth of the hash table as the maximum number of allowable refinement levels ( $l_{\max}$ ) increases. The worst case spatial complexity for our 2-D problem goes as  $O(4^{l_{\max}} n)$ . However, this can be resolved by “local” hashing, which partially constructs the hash table in local regions of the mesh, or by “iterative” hashing, which constructs the hash table iteratively by level of refinement, thus exploiting the fact that the entire hash table need not be constructed at once for our numerical methods. Implementing these hashes is a promising avenue of research.

The impressive speed-ups prompted further research into the uses of hashing for numerical methods on GPUs and CPUs. Robey *et al.* [7] extended beyond neighbor calculation to consider sorting, remapping, and table look-up as well. An important point discussed there is that hash-based algorithms benefit from a two-fold performance boost, seeing speed-up both in algorithm design and in the ease of parallelization without recourse to scan operations (at least in the neighbor calculation algorithm).

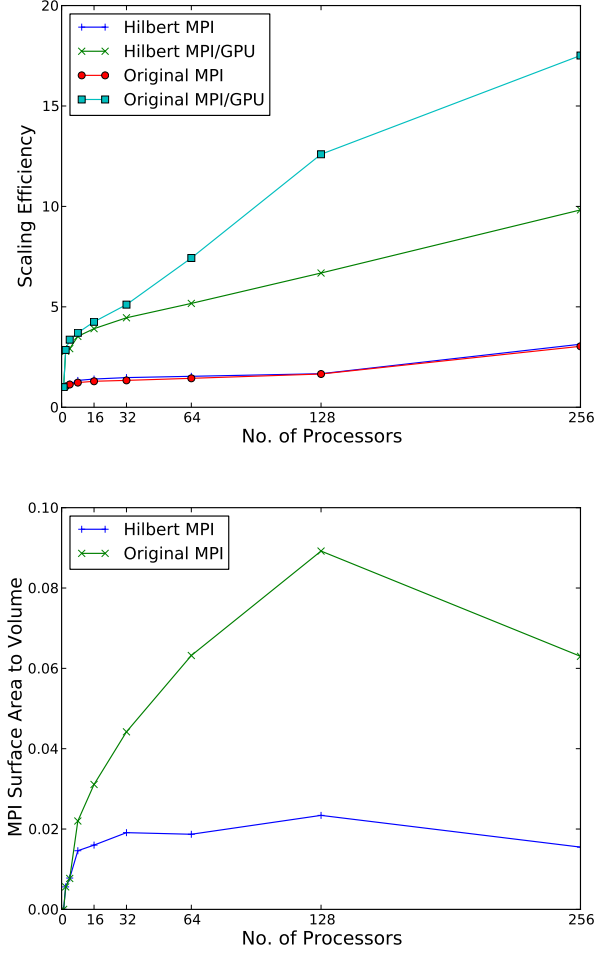


Fig. 11. The upper graph shows weak scaling for a problem run on a  $512 \times 512 \times N$  coarse mesh ( $N$  equal to the number of nodes) with a maximum of 2 levels of refinement for 500 iterations. The lower graph shows the MPI partition measure as a function of the number of nodes for the problem run.

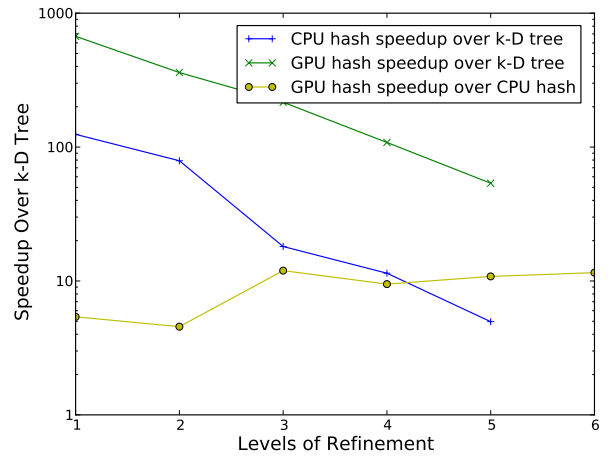


Fig. 12. The speed-up over a  $k$ -D tree for the hash table implementation on both the CPU and GPU.



## V. CONCLUSION

We've shown that it is viable to implement cell-based adaptive mesh refinement on clusters of general purpose graphics processing units. We've shown that finite differencing in this scheme requires careful consideration of the fluxes in order to ensure the conservation of the state variables. We've also stressed the necessity of thinking locally. We took data ordering as the free parameter to establish a partitioning of the mesh which is the most efficient, and then we addressed the issues it created. Namely, the nontrivial neighbor accessing was resolved via a hash-based method of neighbor indexing. Ultimately, we've shown that GPGPUs can be used for intense numerical calculations while making memory and not FLOPs the primary focus in algorithm architecture. The order of magnitude overall speed-up is indicative of a methodology worth the effort.

## ACKNOWLEDGEMENTS

We would like to thank Dennis Trujillo for his collaboration in developing the concepts for the CLAMR code, Rachel Robey for the real-time OpenGL graphics code, H. C. Edwards of Sandia National Laboratories for the global Hilbert space-filling curve code, and Richard Barrett for the sparse MPI communication package. Also helpful were discussions with the Applied Math Department at the University of Washington, colleagues at Lawrence Livermore National Lab, the University of California at Davis, and fellow colleagues at LANL, that helped refine our work as it was developing. The lead author would also like to thank James Quirk for invaluable discussions on computational classics.

## APPENDIX: OPENCL NOMENCLATURE

Although intended to be device-independent, computer architecture naturally dictates that OpenCL distinguish between the *host* and the *device*. The host sets up *kernels*, or GPU-based functions, which are executed on the device. The device should be highly multithreaded, and the kernel should be structured so as to take advantage of this data parallelism and the GPU's high memory bandwidth. OpenCL denotes the independent processes or threads as *work-items*, which are themselves organized into clustered *workgroups*. Within a workgroup execution cycle, blocks of code are executed in lock-step, so efficient kernels should be written with the priority of minimizing process branching and random global memory accesses, and privileging data parallelism and local and shared memory accesses (random or vectorized). Latency is hidden by data processing rather than by use of a cache system as on the CPU. The GPU may be thought of as having a programmable cache, local on-chip memory, as highlighted by efficient data-local methods discussed in Section I.

NVIDIA has published several guides to OpenCL programming and best practices, which should be consulted [32], [10].

## REFERENCES

- [1] Khronos Group, *The OpenCL Specification*, version 1.1, revision 44 ed., Jun. 2011.
- [2] G. Shainer, A. Ayoub, P. Lui, and T. Liu, "Raising the speedlimit: New GPU-to-GPU communications model increases cluster efficiency," Jan. 2011.
- [3] N. E. Davis, R. W. Robey, C. R. Ferenbaugh, D. Nicholaeff, and D. P. Trujillo, "Paradigmatic shifts for exascale supercomputing," *The Journal of Supercomputing*, pp. 1–22, 2012. [Online]. Available: <http://dx.doi.org/10.1007/s11227-012-0789-3>
- [4] J. J. Quirk, "AMR\_sol: Design principles and practice," In *29th Computational Fluid Dynamics, VKI Lecture Series, Chapter 5*, 1998.
- [5] M. J. Berger and J. Olinger, "Adaptive mesh refinement for hyperbolic partial differential equations," *Journal of Computational Physics*, vol. 53, no. 3, pp. 484 – 512, 1984.
- [6] M. J. Berger and P. Colella, "Local adaptive mesh refinement for shock hydrodynamics," *Journal of Computational Physics*, vol. 82, pp. 64–84, May 1989.
- [7] R. N. Robey, D. Nicholaeff, and R. W. Robey, "Hash-based algorithms for discretized data," In *Revision, SIAM J. of Sci. Com.*, 2013.
- [8] D. E. Knuth, *The art of computer programming. Vol. 1, Fundamental algorithms*, 3rd ed. Addison-Wesley Pub. Co., Jul. 1997.
- [9] NVIDIA, "OpenCL programming guide for the cuda architecture, version 3.2," Aug. 2010. [Online]. Available: [http://www.nvidia.com/content/cudazone/download/OpenCL/NVIDIA\\_OpenCL\\_ProgrammingGuide.pdf](http://www.nvidia.com/content/cudazone/download/OpenCL/NVIDIA_OpenCL_ProgrammingGuide.pdf)
- [10] —, "OpenCL best practices guide," May 2010. [Online]. Available: [http://www.nvidia.com/content/cudazone/CUDABrowser/downloads/papers/NVIDIA\\_OpenCL\\_BestPracticesGuide.pdf](http://www.nvidia.com/content/cudazone/CUDABrowser/downloads/papers/NVIDIA_OpenCL_BestPracticesGuide.pdf)
- [11] T. Tao, "The shallow water wave equation and tsunami propagation," <http://terrytao.wordpress.com/2011/03/13/>, 2011.
- [12] L. D. Landau and E. M. Lifshitz, *Fluid Mechanics, Second Edition: Volume 6*, 2nd ed., ser. Course of theoretical physics. Butterworth-Heinemann, Jan. 1987.
- [13] M. J. Castro Díaz, S. Ortega Acosta, M. d. I. Asunción, J. M. Mantas, and J. M. Gallardo, "GPU computing for shallow water flow simulation based on finite volume schemes," *Comptes Rendus Mécanique*, vol. 339, no. 2, pp. 165–184, 2011.
- [14] A. R. Brodtkorb, M. L. Sætra, and M. Altinakar, "Efficient shallow water simulations on gpus: Implementation, visualization, verification, and validation," *Computers & Fluids*, vol. 55, no. 0, pp. 1 – 12, 2012. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0045793011003185>
- [15] P. Lax and B. Wendroff, "Systems of conservation laws," *Communications on Pure and Applied Mathematics*, vol. 13, no. 2, pp. 217–237, 1960.
- [16] R. LeVeque, *Finite volume methods for hyperbolic problems*, ser. Cambridge texts in applied mathematics. Cambridge University Press, 2002.
- [17] D. George, "Numerical approximation of the nonlinear shallow water equations with topography and dry beds: a Godunov-type scheme," Master's thesis, 2004.
- [18] S. F. Davis, "A simplified TVD finite difference scheme via artificial viscosity," *SIAM Journal on Scientific and Statistical Computing*, vol. 8, no. 1, pp. 1–18, 1987.
- [19] P. K. Sweby, "High resolution schemes using flux limiters for hyperbolic conservation laws," *SIAM J. Numer. Anal.*, vol. 21, no. 5, pp. 995–1011, 1984.
- [20] H. Yee, "Construction of explicit and implicit symmetric tvd schemes and their applications," *Journal of Computational Physics*, vol. 68, no. 1, pp. 151 – 179, 1987.
- [21] B. Coutinho, D. Sampaio, F. Pereira, and W. Meira, "Performance debugging of gpgpu applications with the divergence map," in *22nd International Symposium on Computer Architecture and High Performance Computing*, Oct. 2010, pp. 33 – 40.
- [22] G. Peano, "Sur une courbe, qui remplit toute une aire plane," *Mathematische Annalen*, vol. 36, pp. 157–160, 1890, 10.1007/BF01199438.
- [23] H. J. Haverkort and F. van Walderveen, "Locality and bounding-box quality of two-dimensional space-filling curves," *CoRR*, vol. abs/0806.4787, 2008.
- [24] G. Jin and J. Mellor-Crummey, "Using space-filling curves for computation reordering," *Proceedings of the Los Alamos Computer Science Institute Sixth Annual Symposium*, 2005.
- [25] D. Hilbert, "Ueber die stetige abbildung einer line auf ein flächenstück," *Mathematische Annalen*, vol. 38, pp. 459–460, 1891, 10.1007/BF01199431.
- [26] G. Morton, *A computer oriented geodetic data base and a new technique in file sequencing*. International Business Machines Co., 1966.
- [27] A. Godiyal, J. Hoberock, M. Garland, and J. Hart, "Rapid multipole graph drawing on the GPU," in *Graph Drawing*, ser. Lecture Notes in

- Computer Science. Springer Berlin / Heidelberg, 2009, vol. 5417, pp. 90–101.
- [28] D. A. Alcantara, A. Sharf, F. Abbasinejad, S. Sengupta, M. Mitzenmacher, J. D. Owens, and N. Amenta, “Real-time parallel hashing on the gpu,” *ACM Transactions on Graphics (Proceedings of ACM SIGGRAPH Asia 2009)*, vol. 28, no. 5, Dec. 2009.
  - [29] H. Ji, F.-S. Lien, and E. Yee, “A new adaptive mesh refinement data structure with an application to detonation,” *Journal of Computational Physics*, vol. 229, no. 23, pp. 8981–8993, 2010.
  - [30] R. W. Robey, J. M. Robey, and R. Aulwes, “In search of numerical consistency in parallel programming,” *Parallel Comput.*, vol. 37, pp. 217–229, April 2011.
  - [31] J. Hungershöfer and J.-M. Wierum, “On the quality of partitions based on space-filling curves,” in *Computational Science - ICCS 2002*, ser. Lecture Notes in Computer Science, P. Sloot, A. Hoekstra, C. Tan, and J. Dongarra, Eds. Springer Berlin Heidelberg, 2002, vol. 2331, pp. 36–45. [Online]. Available: [http://dx.doi.org/10.1007/3-540-47789-6\\_4](http://dx.doi.org/10.1007/3-540-47789-6_4)
  - [32] NVIDIA, “OpenCL programming overview,” Aug 2009. [Online]. Available: [http://www.nvidia.com/content/cudazone/download/OpenCL/NVIDIA\\_OpenCL\\_ProgrammingOverview.pdf](http://www.nvidia.com/content/cudazone/download/OpenCL/NVIDIA_OpenCL_ProgrammingOverview.pdf)