

TECHNICAL REPORT

Uintah: A Scalable Adaptive Framework for Emerging Petascale Platforms

J. Luitjens and M. Berzins

UUSCI-2009-002

Scientific Computing and Imaging Institute
University of Utah
Salt Lake City, UT 84112 USA

July 1, 2009

Abstract:

Uintah is a highly parallel and adaptive mesh multi-physics framework created by the Center for Simulation of Accidental Fires and Explosions. Uintah, which is built upon the Common Component Architecture (CCA), has facilitated the simulation of a wide variety of fluid-structure interaction problems using both structured adaptive meshes to model fluids and particles to model solids. Uintah was originally designed for, and has performed well on, a few thousand cores. The extension of Uintah to use tens of thousands cores has required improvements in memory usage, data structure design, and load balancing algorithms. These improvements have led to improved strong and weak scalability up to 32,768 cores. This paper will describe these improvements, including a memory tracking system and a novel cost estimation algorithm that utilizes forecasting for dynamic load balancing, and demonstrate Uintah's improved scalability on a two material compressible Navier Stokes problem.

Uintah: A Scalable Adaptive Framework for Emerging Petascale Platforms.

Justin Luitjens
School of Computing
University Of Utah
Salt Lake City, Utah 84112
luitjens@cs.utah.edu

Martin Berzins
School of Computing
University Of Utah
Salt Lake City, Utah 84112
mb@sci.utah.edu

ABSTRACT

Uintah is a highly parallel and adaptive mesh multi-physics framework created by the Center for Simulation of Accidental Fires and Explosions. Uintah, which is built upon the Common Component Architecture (CCA), has facilitated the simulation of a wide variety of fluid-structure interaction problems using both structured adaptive meshes to model fluids and particles to model solids. Uintah was originally designed for, and has performed well on, a few thousand cores. The extension of Uintah to use tens of thousands cores has required improvements in memory usage, data structure design, and load balancing algorithms. These improvements have led to improved strong and weak scalability up to 32,768 cores. This paper will describe these improvements, including a memory tracking system and a novel cost estimation algorithm that utilizes forecasting for dynamic load balancing, and demonstrate Uintah's improved scalability on a two material compressible Navier Stokes problem.

1. INTRODUCTION

The University of Utah Center for the Simulation of Accidental Fires and Explosions (C-SAFE) [1] is a Department of Energy ASC center that focuses on providing state-of-the-art, science-based tools for the numerical simulation of accidental fires and explosions. The primary objective of C-SAFE has been to provide a software system in which fundamental chemistry and engineering physics are fully coupled with nonlinear solvers and visualization tools, thereby integrating expertise from a wide variety of disciplines. The creation of Uintah has furthered C-SAFE's understanding of fires, explosions, and other problems involving complex fluid-structure interactions.

For example, on August 11, 2005 a truck carrying 35,500 pounds of explosives down Utah's Spanish Fork Canyon overturned and caught fire. Within minutes the truck detonated with a force much larger than expected leaving behind a 70 foot crater. Fortunately no one was hurt. Why did a detonation occur as opposed to a deflagration, which is several

orders of magnitude less violent? Could the packing of the individual explosive charges influence the propagation of the combustion wave, or the amount of energy released? These are the types of questions that C-SAFE is addressing and hopes to further address through the use of future petascale simulations. Large-scale simulations have allowed C-SAFE to further the understanding of explosions by providing the ability to look more closely at the underlying physical phenomena than is possible through experimental tests.

The target simulation scenario for C-SAFE is a small cylindrical steel container filled with plastic bonded explosive (PBX-9501) subjected to convective and radiative heat fluxes from a fire which heats the container and the PBX. After some amount of time the critical temperature in the PBX is reached and the explosive begins to rapidly decompose into a gas. The solid-to-gas reaction pressurizes the interior of the steel container causing the shell to rapidly expand and eventually rupture. The gaseous products of reaction form a blast wave that expands outward along with pieces of the container and the unreacted PBX. This scenario along with images from a Uintah simulation can be seen in Figure 1.

Simulating this problem requires expertise from a wide variety of disciplines including combustion, structural mechanics, and fluid dynamics. In addition, such a problem requires a large amount of processing power necessitating the need for both adaptive mesh refinement (AMR) [2] and parallelism. AMR focuses the computational resources where needed by adding refinement in areas where rapidly evolving physical processes are occurring. For example, in the case of the exploding container mentioned above; the container, the explosive, and the pressure wave all need to be highly resolved, where as the surrounding atmosphere has a lower resolution requirement. Even with AMR, the processing requirements for such a problem are still large necessitating the use of parallelism.

The need for parallelism, AMR, and a wide variety of physics has led to the development of the Uintah Computational Framework [1, 3, 4]. Uintah, which was developed by C-SAFE under the direction of Steven Parker, provides a large degree of encapsulation that allows scientists to focus on their area of expertise without fully understanding complexities outside of their domain.

In preparation for petascale architectures and simulations, the performance of frameworks like Uintah must be ana-

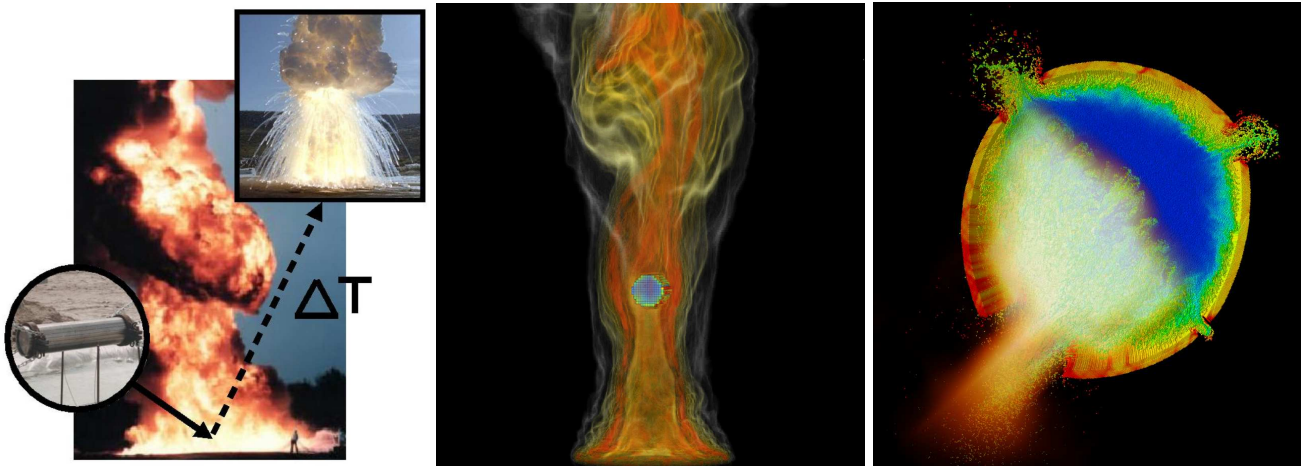


Figure 1: Images of C-SAFE’s target problem: a container PBX is heated up by a pool fire until the explosive ignites, pressurizing the container, and causing the container to rupture. As the container ruptures, a blast wave containing gaseous products, parts of the container, and unreacted PBX expands outward.

lyzed and optimized. Poor performance in any portion of the framework can have a significant impact on the overall performance. Achieving a high degree of scalability for AMR based simulations is challenging due to poor scalability associated with the changing grid. With AMR, whenever the grid changes, a number of operations must be performed. For example, the new grid must be created, work must be load balanced and migrated to the owning cores, and the communication schedule must be created. In the past, optimizations to these operations have led to improvements in scalability within general purpose AMR frameworks [5, 6, 7]. Previously application specific AMR codes have been shown to scale up to 60,000 cores [8, 9]. The challenge is to see if a general purpose framework such as Uintah is capable of scaling to such numbers of cores. We address this challenge by: describing Uintah and its novel approach to parallelism, describing a tool used to identify inefficiencies in memory usage and data structures, and presenting a new method to estimate costs used in load balancing, which together have led to substantial improvements in Uintah’s scalability.

2. UINTAH

The Uintah computational framework, is a set of parallel software components and libraries built upon the DOE Common Component Architecture (CCA) that facilitate the solution of partial differential equations (PDEs) on structured AMR grids. Uintah is a sophisticated framework that can integrate multiple simulation components, analyze the dependencies and communication patterns between them, and efficiently execute the resulting multi-physics simulation. Uintah employs an abstract task graph representation to describe computation and communication [3, 4]. Through this mechanism, Uintah components delegate decisions about parallelism to a framework component, which determines communication patterns and characterizes the computational workloads needed for global resource optimization. This allows parallelism to be integrated between multiple components while maintaining overall scalability. Uintah also analyzes the structure of the computation and automatically enables load balancing, data communication, parallel

I/O, checkpointing and restarting capabilities.

One of the primary strengths of Uintah is that application designers can develop large-scale parallel SAMR simulations with little understanding of the underlying parallelism. To do this the designers must specify their algorithm as a series of serial tasks that run on a hexahedral mesh patch. Each task specifies the computation to be performed for a single time step and the related variable dependencies. Variable dependencies state what variables the task requires for the computation (along with the stencil width) and what variables the task modifies or computes. Using these dependencies, Uintah creates a directed acyclic task graph that specifies the task execution order and the required communication for the simulation. This design shields developers from the parallelism while allowing Uintah to utilize highly sophisticated communication patterns including a large amount of asynchronous communication and message coalescing. By using these advanced communication techniques Uintah is able to hide the cost of some of the communication by overlapping it with computation. As we move to petascale architectures advanced communication techniques like asynchronous communication will, most likely, be increasingly necessary.

Uintah achieves parallelism by dividing the grid into hexahedral mesh patches, which are uniquely assigned to cores. Each core executes the tasks on its assigned patches achieving a domain-based parallelism. When a task requires a variable with a non-zero stencil width, communication between neighboring patches is required before the task can execute. Using the task graph and the core assignments for each patch, Uintah determines the communication necessary and schedules communication and computation at the appropriate times. All communication, including intra-level (within a single level), inter-level (between AMR levels), and data migration after load balancing is included in this schedule. The schedule is then executed repeatedly with each execution corresponding to a single time step of the simulation. In static grid computations, this schedule is created once and

reused for the entire simulation. However, in SAMR computations the schedule must be recreated whenever the patch set changes (regridding) or whenever the patch assignments change (load balancing). The framework also utilizes parallel I/O to store simulation data for use in checkpointing, restarting, and visualization within VisIt [10].

Uintah was originally designed for a few thousand cores and has been used regularly for simulations with up to 2,000 cores. However, scalability at larger numbers of cores was problematic because the memory utilization, data structures, and load balancing did not scale well on 4,000 or more cores. In particular, memory utilization was an issue due to data structures that consumed memory on the order of the number of cores or the number of patches. As the number of cores or patches increased the size of these data structures would also increase eventually exceeding the available resources. Recently these inefficiencies were resolved through the creation of a tool that aides in the tracking memory allocations over time. These improvements are described in Section 3.

The component design has allowed Uintah to excel as a research platform. Components can be swapped in and out, allowing them to be developed and tested within the entire framework, without affecting other components. This has led to a highly flexible simulation package which has been able to simulate a wide variety of problems including shape charges, stage-separation in rockets, the biomechanics of microvessels [11], the properties of foam under large deformation [12], and the evolution of large pool fires caused by transportation accidents [13], in addition to the exploding container scenario described in Section 1.

Uintah currently contains three main simulation algorithms, or components, that are capable of using AMR: i) the ICE compressible multi-material CFD formulation [14, 15, 16], ii) the particle-based Material Point Method (MPM) [17] for structural mechanics, and iii) the combined fluid-structure interaction algorithm MPMICE [18]. In addition, Uintah integrates numerous sub-components including equations of state, constitutive models, and reaction models.

ICE is a “multi-material” CFD algorithm that was developed by Kashiwa and others at LANL [14, 15, 16]. This technique can be used in both incompressible and compressible flow regimes, which is necessary when modeling fires and explosions. The cell centered, finite volume solution technique is convenient in that a single control volume is used for all materials. Conserving mass, momentum and energy, and the exchange of these quantities between the materials is simplified by use of a common control volume.

The *Material Point Method* is a particle method that is used to evolve the equations of motion for the solid materials. MPM is a powerful technique for computational solid mechanics, and has found favor in many applications involving complex geometries [11], large deformations [12], and fracture [19]. Originally described by Sulsky, et al., [20], MPM is an extension to solid mechanics of FLIP [21, 22], which is a particle-in-cell (PIC) method for fluid flow simulation [23].

3. MEMORY IMPROVEMENTS

Previously, scalability of Uintah to larger numbers of cores was problematic due to memory inefficiencies. Portions of the framework used data structures that had memory complexity on the order of the number of patches or cores. As the problem size or number of cores increased the memory requirement of these data structures would also increase.

Uintah has been in development for over ten years and now contains around a half a million lines of code. The development process has involved many different people, many of which are no longer working on Uintah. Uintah was initially designed for upwards of a thousand cores. The data structures and algorithms used in Uintah worked well for a few thousand cores but became inefficient when moving to tens of thousands of cores. Inefficiencies like these will be common in many codes attempting to move onto petascale platforms. Identifying these inefficiencies in large legacy codes like Uintah is a daunting task. Traditional tools, like TAU, have a large memory footprint preventing its usage on large problems in Uintah. This has led to the development of MallocTrace by one of the authors [24].

MallocTrace is a low-memory-overhead tool for tracking memory usage within an application. This tool logs memory allocations in C++ programs through a series of macros and library hooks. The logs contain information including file names and line numbers which are useful for tracking memory usage. The tool also provides a basic mechanism to parse the logs and provides a summary of memory usage at any given time in the simulation. This allows a user to see where memory is allocated at any point in the simulation. The low memory overhead allows the tool to be used with a large number of cores on programs like Uintah.

MallocTrace allowed us to rapidly identify and eliminate inefficiencies in memory usage at large numbers of cores. The effect of eliminating these inefficiencies can be seen in Figure 2. This graph shows an algorithmic decrease in memory usage. Prior to the optimizations the memory usage would sometimes increase with the number of cores. The same increase was not seen after the optimizations.

Eliminating these memory inefficiencies not only lowered the size of Uintah’s memory footprint, but also provided significant improvements in runtime. Figure 3 show the decrease in runtime due to these optimizations of an ICE simulation of an expanding 3D blast wave using AMR. Prior to the optimizations the runtime rose sharply after 2048 cores. After the optimizations the same increase was not seen and Uintah continued to scale to 4096 cores.

Through the use of the MallocTrace tool we quickly identified and eliminated memory inefficiencies leading to significant increases in performance. It took less than a week from the deployment of this tool to achieve the results shown here. Inefficiencies like those found in Uintah are not uncommon in large legacy codes. Undoubtedly similar issues will exist when moving to hundreds of thousands of cores. Determining which portions of Uintah need to be redesigned is a challenging task that is made easier with tools that are capable of running on large-scale problems like MallocTrace. For more information on MallocTrace please see [24].

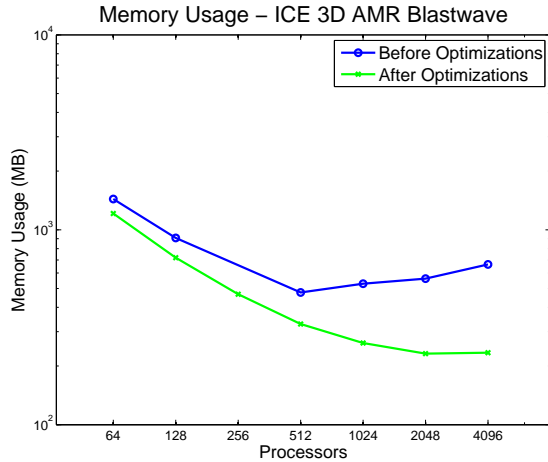


Figure 2: A comparison of the memory usage of Uintah before and after the elimination of inefficiencies identified with MallocTrace.

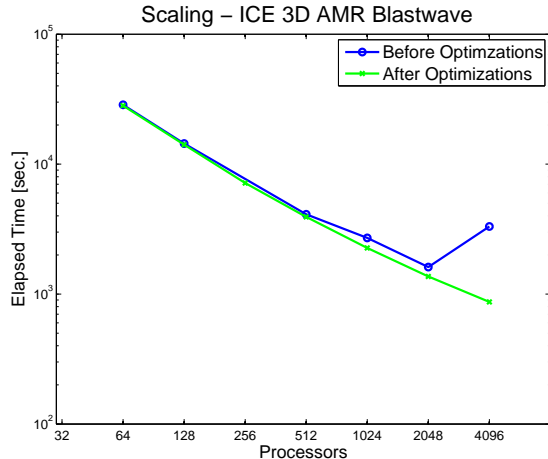


Figure 3: A comparison of the scalability of Uintah before and after the elimination of inefficiencies identified with MallocTrace.

4. DYNAMIC LOAD BALANCING

The variety of available simulation components within Uintah requires a sophisticated load balancer that is flexible enough to handle all of Uintah simulations. Dynamic load balancing can be described as the minimization of three competing costs: The cost of load imbalance, the cost of communication, and the cost to generate the load distribution.

A load imbalance occurs when one or more cores are assigned more work than other cores. A large load imbalance will cause cores to wait for other cores to finish their computation leading to poor utilization of system resources.

In addition, too much communication can also cause performance issues. Communication across the network is slow relative to the time for computation and can easily dominate the time to reach a solution. In many simulations, communication is predominantly local, meaning that only a small area around each patch must be communicated from physically neighboring patches. By clustering neighboring patches together the framework can greatly reduce the necessary communication and significantly affect the overall runtime.

Finally, with AMR methods the workload changes as the mesh changes. In addition, with particle methods the workload can change on each time step as particles move throughout the domain. This can cause load balancing to occur often, making it important that the time to generate the patch distribution is small relative to the overall computation. If a slow load balancing algorithm is used and load balancing occurs often, the time to load balance can dominate the overall runtime. In this case, it may be preferable to use a faster load balancing algorithm that produces more load imbalance.

The need for fast and effective load balancing techniques has led to the development of widely used load balancing applications like Metis [25], Jostle [26], and Zoltan [27, 28]. Uintah has recently added support for the Zoltan load balancing package, providing easy access to a number of algorithms. In addition, Uintah can use its own highly parallel load balancing algorithm [29] that utilizes space-filling curves, which has been shown to be better than Zoltan’s space-filling curve load balancer within Uintah [30].

To balance the computation effectively load balancers need an estimate of the cost (execution time) of the computation. A poor estimate of the cost will lead to a decrease in load balance. Thus it is important that this estimate be accurate. One method to estimate these costs is to use algorithmic cost models.

4.1 Algorithmic Cost Models

Algorithm cost models (ACM) attempt to model the underlying algorithms. For example, the ICE algorithm is a cell-based algorithm that performs a constant amount of work per cell and as such the cost is proportional to the number of cells. Equation (1) below, describes an accurate ACM for ICE, where C_p is the cost of a patch, N_c is the number of cells in that patch, and c_1 is the constant time execution time the ICE algorithm on a single cell.

$$C_p = c_1 N_c \quad (1)$$

In addition to performing cell-based computations, MPMICE also has particle-based computations in regions where solid materials exist. Equation 2 below describes a possible ACM for MPMICE, where N_p is the number of particles within the patch and c_2 is the constant execution time on a particle.

$$C_p = c_1 N_c + c_2 N_p \quad (2)$$

This model is not as accurate as the ICE model because the work performed by MPMICE is not constant per particle or cell. In MPMICE, during the equilibration pressure solve, the simulation performs a local iterative solve on a per-cell basis [18]. This solve may converge at different rates throughout the domain depending on the underlying physics. Capturing such behavior in an ACM is a challenging task. Furthermore the motion of particles between cells can cause the workload per core to change at every time step and not only when regridding occurs, as is the case in AMR simulations.

In addition to developing these models, estimates for the constants must be determined on a per-problem basis. The constants can vary greatly depending on the underlying physical processes. To make matters worse these constants can also vary according to system architectures, compilers, and compiler options. In order to achieve an effective load balance, these constants must be proportionally accurate. For models with a single constant, like the model used for ICE, estimating the constant is trivial. However, the difficulty in estimating these constants increases significantly with the number of constants in the model. Maintaining an accurate list of these constants for each possible problem, architecture, and compiler combination is not feasible, thus placing the challenge of estimating these constants on the user.

4.2 Forecasting Cost Model

Since developing an accurate algorithmic cost model is challenging, we have added an alternative approach to Uintah which utilizes a forecasting cost model (FCM) to predict the cost of each patch based on time series. During task execution, the time to complete each task on a region of the domain is recorded and used to update a simple forecasting model. That model is then used to predict the execution time on that region in the future. This provides a mechanism to accurately predict the cost of each patch and eliminates the need to estimate constants for an ACM. Uintah uses simple exponential smoothing as its forecasting model [31]. The model is as follows:

$$W_{r,t+1} = \alpha E_{r,t} + (1 - \alpha)W_{r,t}, \quad (3)$$

where $W_{r,t}$ is the predicted cost at time step t on region r , $E_{r,t}$ is the actual execution time at time step t on region r , and α is a weighting factor in the range of $[0,1]$ which represents the rate of decay on past data. This method

can also be viewed as a weighted moving average where the weight on past observations decreases exponentially [31]. A smaller value for α causes the algorithm to put more weight on recent observations causing the forecast to respond more quickly to changes in the actual value but also causes the forecast to become more susceptible to noise. A larger value for α will cause that data to be smoother eliminating noise but also causes the forecast to react more slowly to changes in the actual value. α can be defined in terms of the size of a moving average window using the following equation:

$$\alpha = \frac{2.0}{T + 1}, \quad (4)$$

where T is the number of time steps that will contain 99.9% of the total weight in the weighted average [31]. Uintah uses a default value of 10 for T .

On the first time step of the simulation $W_{r,t}$ is unavailable, requiring an estimation of the initial value. For the initial time step Uintah load balances using the algorithmic cost models described above. The initial measurements are then used to set the initial value by setting $W_{r,0} = E_{r,0}$.

A different initialization approach is used when new regions of refinement are created during the regridding process. During this process refinement may be added in regions where it previously did not exist. When these regions are created, $W_{r,t}$ must be estimated. Using an ACM would likely produce a poor estimate that would not be proportionally accurate to the forecasted values elsewhere in the domain. In order to estimate the cost while maintaining proportional accuracy, Uintah sets $W_{r,t}$ for the new regions equal to the average value of $W_{r,t}$ for all regions. This ensures that the initial value for the new region is at least close to the actual value which also ensures that the estimation will be accurate within a few time steps and that load imbalance caused by this estimation will be limited.

To allow for changing patch sets forecasting is performed on a per-region basis instead of a per-patch basis. The difference between regions and patches is shown in Figure 4. Regions are constant-sized portions of the domain that are contained within a single patch. Patches on the other hand are variable-sized portions of the domain that may contain many regions.

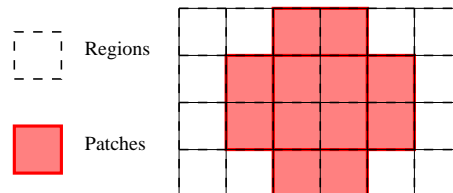


Figure 4: The difference between regions and patches. Patches are composed of one or more fixed size blocks referred to as regions.

By forecasting on a per-region basis, the patch set can change without needing to interpolate forecasting data between the

changing patch sets. This necessitates mechanisms to interpolate the data between regions and patches. Since regions are completely contained within patches these mechanisms are straight forward to describe and implement. The measured cost for each region is equal to the cost of the patch, times the proportion of the patch that the region encompasses, as described in the following equation:

$$E_{r,t} = E_{p,t} \frac{V_r}{V_p}, \quad (5)$$

where t is the current time step, $E_{r,t}$ is the measured computation time for region r , $E_{p,t}$ is the measured execution time for patch p , V_p is the volume of patch p , and V_r is the volume of region r . In addition, $W_{p,t+1}$ can be defined as the sum of the weights of all regions contained in patch p :

$$W_{p,t+1} = \sum_r^{r \in p} (W_{r,t+1}). \quad (6)$$

Uintah stores the forecasting data while minimizing both storage and communication. The forecasting data is stored locally on each core. When a core executes a task on a patch, it adds the contribution to its local forecast data using Equation (5). If a region was owned by a different core in the past, then local forecast data will exist on multiple cores but each core will only update its local data. At the end of each time step the simulation finalizes the forecast data by applying Equation (3). Updating the forecast data each time step is a local operation which does not require any communication. However, communication is required when load balancing occurs. During load balancing, each core must know the cost of each patch. This is done by applying Equation (6) locally and then performing a `MPLAllreduce` to get the global sum.

In order to keep the data structures for forecasting as small as possible, contributions are stored in a Standard Template Library (STL) map, which is a sparse data structure. This causes the storage per core to be proportional to the number of patches per core. In addition, when the contribution for a region becomes too small it is deleted from the map. When a core has not updated a region in its map for over T time steps, the contributing weight for that region is less than 0.1% of the total weight, at this point we consider the weight to be insignificant and delete it from the map. This prevents the size of the maps from slowly increasing over time.

4.3 Forecasting Results

The effect of forecasting on Uintah’s runtime was tested using two different simulation components. The first test used the ICE, multi-material algorithm with explicit time stepping to simulate the transport of two fluids with a prescribed initial velocity. For this problem the conservation of mass, momentum, and energy equations are solved for two inviscid fluids. The fluids exchange momentum and heat through the exchange terms in the governing equations. This problem exercises all of main features of ICE and amounts to

solving eight P.D.E’s, along with two point-wise solves, and one iterative solve for more information see [18].

The second simulation was C-SAFE’s target problem using the MPMICE algorithm with explicit time stepping seen in Figure 1. In this problem, a steel container filled with an explosive material is suspended over a fire. As the simulation progresses, the explosive heats up and ignites causing the container to rupture resulting in a violent explosion. This is a complex problem whose computational cost is difficult to predict. As the container ruptures, the performance characteristics of the problem rapidly change as pressurized gasses and explosive materials move across the domain. These problems were selected because of the complexity of the relevant physics.

The computational cost of the ICE problem is predictable and developing an accurate ACM is straight forward. In contrast, the computational cost of the MPMICE problem is difficult to predict, hindering the creation of an accurate ACM. For the ICE simulation, Equation (1) above was used for the ACM with $c_1 = 1$. For the MPMICE simulation, Equation (2) above was used for the ACM with $c_1 = 1$ and $c_2 = 1.25$. While these values are not representative of actual machine constants, they are proportionally accurate, which is sufficient for the load balancing process.

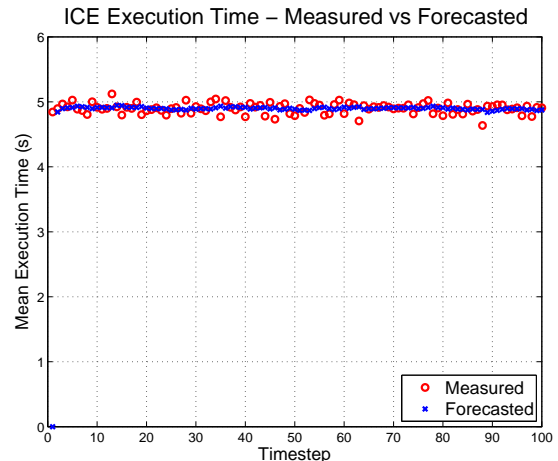


Figure 5: A comparison of the measured and forecasted times spent executing tasks at each time step for ICE.

The measured and forecasted execution times for ICE are shown in Figure 5 and for MPMICE are shown in Figure 6. The measured time for both problems fluctuates due to system noise causing error in the forecast. The forecast error was on average 7% for the ICE problem and 5% for the MPMICE problem. In both of these tests forecasting does an effective job predicting computation time.

Figure 7 shows the difference in load imbalance for the ICE simulation using a FCM versus an ACM. The load imbalance varies between 3% and 10% with an average imbalance of 5.3% using the FCM and 6% using the ACM. In this case the difference in runtime was marginal. This shows that Uintah’s performance using a FCM is similar to performance

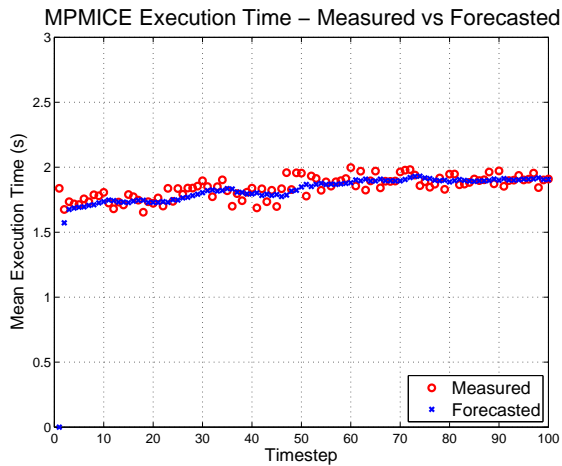


Figure 6: A comparison of the measured and forecasted times spent executing tasks at each time step for MPMICE.

using an accurate ACM.

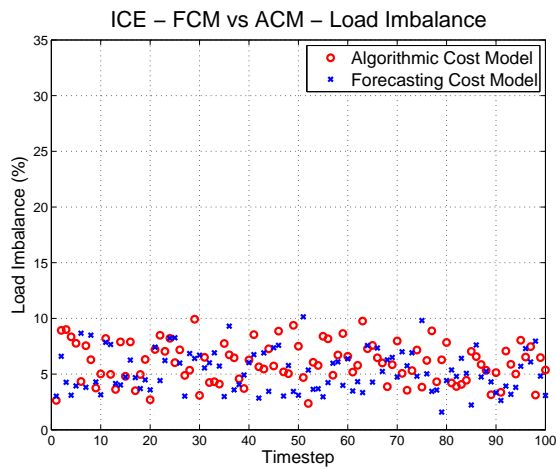


Figure 7: A comparison of the load imbalance when using an ACM versus a FCM for ICE.

The load imbalance for the MPMICE simulation can be seen in Figure 8. When using an ACM, the load imbalance varies between 13%-35% with an average imbalance of 20%. This variance is due to rapid changes in the performance characteristics that are not captured by the current model. At the same time the load imbalance when forecasting was relatively constant with an average imbalance of 4%, which is consistent with the fluctuations we saw in the measurements in Figure 6.

The improved load balance led to a substantial increase in performance seen in Figure 9. When forecasting, the MPMICE simulation was approximately 15% faster than when using the ACM.

These results show that forecasting can produce accurate cost estimations that are at least as effective as an accurate

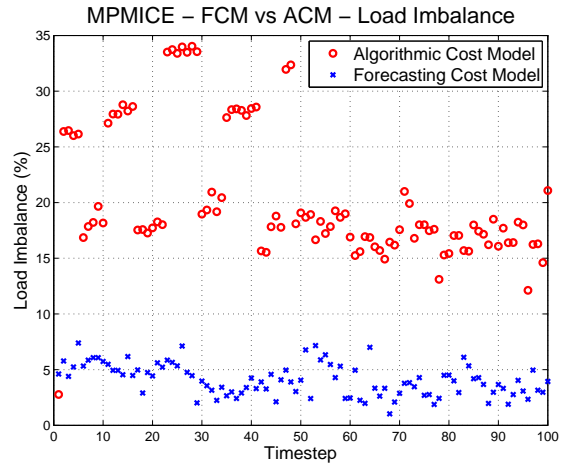


Figure 8: A comparison of the load imbalance when using an ACM versus a FCM for MPMICE.

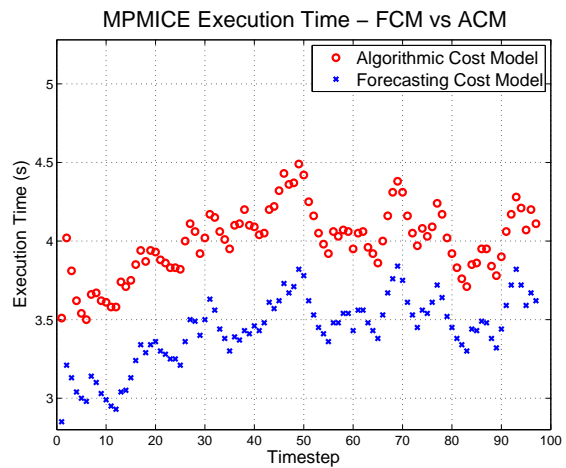


Figure 9: A comparison of the time spent in task execution in MPMICE using an ACM versus a FCM.

ACM. In addition, forecasting is able to predict complex interactions that may be difficult to capture in an ACM leading to improved cost estimations and reduced runtimes. Higher order forecasting methods described in [31] have also been used but provided no benefit over the first order forecasting methods described above. In the future, it may be worthwhile to use a Kalman filter [32] or other advanced forecasting methods.

5. SCALABILITY OF AMR IN UINTAH

The scalability of AMR in Uintah was tested in both the weak and strong sense using the ICE problem described in section 4.3 on Ranger¹. For **weak scaling** the problem size per core is held constant as the number of cores increases and in **strong scaling** the total problem size is held constant while the number of cores increases. For each scaling test the total amount of time to complete twenty time steps of the simulation was recorded. This problem contained three mesh levels with each level being a factor of four more refined than the coarser level. Patches were uniformly sized with 16^3 cells in each patch. Regridding and load balancing occurred for each simulation as needed and occurred on average every 8 time steps. Scalability was tested on five problems with each problem being a factor of four larger than the previous problem. The smallest problem contained 1.7 million cells and the largest problem contained 435 million cells.

Figure 10 shows a break down of Uintah’s strong scalability. In this figure the narrow bar represents the maximum time across all cores and the wide bar represents the average time, with the difference between those two bars representing the load imbalance. The black line represents the total time, which is approximately equal to the sum of the average times. Task execution includes the time spent executing each task without the cost of communication, the global communication includes time within collective operations like MPIAllreduce, the wait time includes the time the simulation spent within MPIWait, and all other times including regridding and load balancing are included in the other time category. This figure shows scalability up to 16,384 cores with scalability tailing off slightly at the last data point. This decrease in scalability is due to an increase in the difference between the maximum and the average execution times (load imbalance) which led to an increase in task wait time (synchronization). The increase is due to the limited number of patches per core (1-2 per core). When the work per core is low, the ability of a load balancer to balance the computation is limited.

A breakdown of Uintah’s weak scaling can be seen in Figure 11. This figure shows that the weak scaling is nearly ideal until the last data point. At the last data point the imbalance on the execution increases, causing the wait time and global communication to also increase (synchronization). The source of the increase in load imbalance is currently being investigated.

Figures 10 and 11 also show an area where Uintah’s performance can be improved. These graphs show that the waiting time is a significant portion of the run time. The wait time

¹Ranger is a supercomputer located at the University of Texas with 62,976 cores. More information is available at <http://www.tacc.utexas.edu/resources/hpcsystems/>.

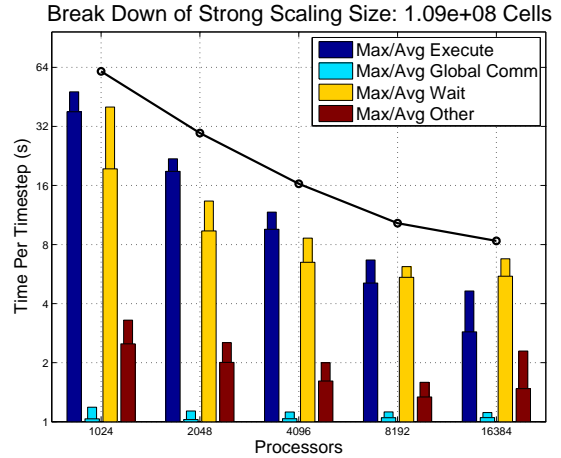


Figure 10: A break down of the strong scaling in Uintah. The thick bar is the average time per core and the narrow bar is the maximum time across all cores.

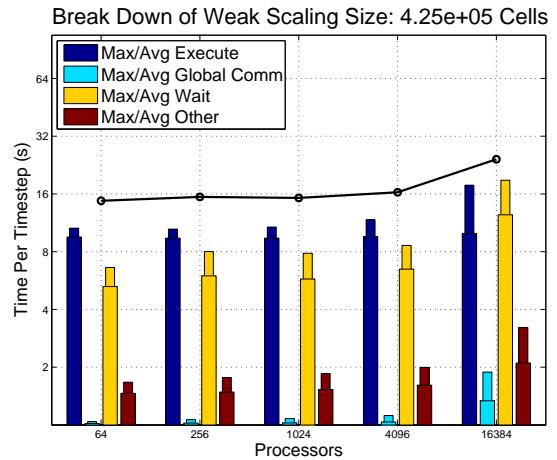


Figure 11: A break down of the weak scaling in Uintah. The thick bar is the average time per core and the narrow bar is the maximum time across all cores.

is defined as time that the simulation is waiting for communication to finish prior to executing a task. The wait time, which can be viewed as point to point communication time, is a combination of time spent waiting for communication to be sent (time for synchronization) and waiting for the communication to arrive (time for communication).

Currently Uintah uses a fixed-order scheduler which pre-determines the order of execution and communication for all tasks. This means that the simulation must wait on the next task to finish communication even when other tasks are ready to execute. Although good scalability has been achieved, experiments have shown that the overall performance could be improved through the use of a dynamic-order schedule that would dynamically schedule tasks when they are ready to execute. This would help decrease both the synchronization and communication costs that are present using the current scheduler. This work is currently being undertaken.

The comprehensive weak and strong scaling up to 32,768 cores can be seen in Figure 12. The elimination of inefficiencies within Uintah and improvements to the load balancer have led to marked improvements to Uintah’s scalability. Strong scaling occurred for every problem size, with a slight tailing off at the last data point due to load imbalance associated with the low amount of work per core. Near ideal weak scaling occurred for the first four data points for all 5 problem sizes, with a sharp uptick occurring at the last data point due to load imbalance that was discussed earlier. This figure shows that Uintah scales in both the weak and strong sense across a large range of problem sizes. It also shows that scalability at larger numbers of cores appears possible. We will continue testing Uintah on larger numbers of cores as they become available. We are currently waiting to test the scalability on all of Ranger’s 62,976 cores.

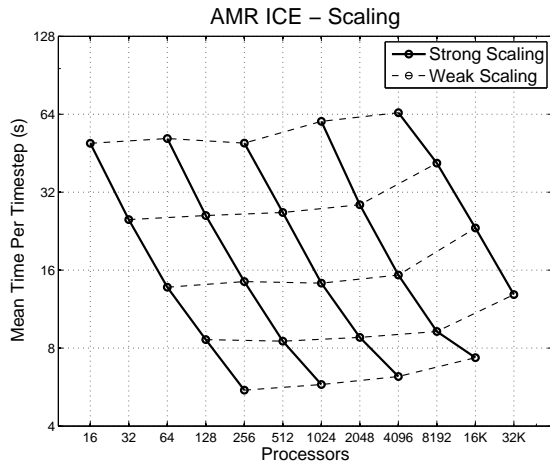


Figure 12: The strong and weak scalability up to 32,768 cores of AMR in Uintah using ICE.

6. CONCLUSIONS AND FUTURE WORK

The primary strength of Uintah is that simulation designers can develop large-scale parallel SAMR simulations with little understanding of the underlying parallelism. This allows for rapid development of large-scale simulations for a

wide variety of problems using Uintah with features like, automated load balancing, parallel I/O, and checkpointing. Uintah’s component design allows for the use of sophisticated algorithms without burdening users by complicating the individual components.

In preparation for emerging petascale architectures, it is important that the performance of frameworks like Uintah are analyzed for inefficiencies. Poor performance in any portion of the framework can hinder performance for the entire simulation. Because of this we have placed substantial effort into identifying, analyzing, and eliminating inefficiencies within Uintah. By analyzing Uintah’s memory usage at large numbers of cores with MallocTrace, we were able to rapidly identify and eliminate multiple inefficiencies. This led to a substantial decrease in memory usage and corresponding increase in performance.

The effect of load balance on the overall performance of a simulation is substantial. A poor load balance will cause poor utilization of system resources preventing scalability. Because of this it is essential that we use effective load balancing algorithms. However, these algorithms will only be effective if the cost estimates provided to them are accurate. Poor cost estimates will cause a poor load balance regardless of what load balancing algorithm is used. While algorithmic cost models can be used to produce these estimates, they are often prone to large error and require the user to estimate model constants. We have shown that an alternative method for estimating these costs eliminates the constants while still providing an accurate estimate that is at least as effective as, and in many cases better than, algorithmic cost models. In the future, more sophisticated forecasting methods could perhaps be used to further improve estimates allowing for a greater utilization of system resources.

Improvements to the memory utilization and load balancer have led to significant improvements in both performance and scalability. We have been able to show both strong and weak scalability of an AMR simulation within a general purpose framework up to 32,768 cores on Ranger and expect even greater scalability when larger machines are made available.

7. ACKNOWLEDGEMENTS

This work was supported by the University of Utah’s Center for the Simulation of Accidental Fires and Explosions (C-SAFE) and funded by both the Department of Energy, under subcontract No. B524196 and the National Science Foundation under subcontract No. OCI0721659. We would like to thank TACC (in particular Karl Schulz) along with Jim Guilkey, Todd Harman, Qingyu Meng, Tom Henderson, and the rest of the C-SAFE team.

8. REFERENCES

- [1] J. D. de St. Germain, S. G. Parker, J. McCorquodale, and C. R. Johnson, “Uintah: A massively parallel problem solving environment,” in *HPDC*, 2000, pp. 33–42.
- [2] M. Berger and P. Colella, “Local adaptive mesh refinement for shock hydrodynamics.” *Journal of Computational Physics*, vol. 82, pp. 64–84, 1989.

- [3] S. G. Parker, J. Guilkey, and T. Harman, "A component-based parallel infrastructure for the simulation of fluid structure interaction," *Engineering with Computers*, vol. 22, no. 3-4, pp. 277–292, 2006.
- [4] S. G. Parker, "A component-based architecture for parallel multi-physics pde simulation," *Future Gener. Comput. Syst.*, vol. 22, no. 1, pp. 204–216, 2006.
- [5] J. Luitjens, B. Worthen, M. Berzins, and T. Henderson, *Petascale Computing Algorithms and Applications*. Chapman and Hall/CRC, 2007, ch. Scalable parallel amr for the uintah multiphysics code.
- [6] A. M. Wissink, R. D. Hornung, S. R. Kohn, S. S. Smith, and N. Elliott, "Large scale parallel structured amr calculations using the samrai framework," in *Supercomputing '01: Proceedings of the 2001 ACM/IEEE conference on Supercomputing (CDROM)*. New York, NY, USA: ACM Press, 2001, pp. 6–6.
- [7] A. M. Wissink, D. Hysom, and R. D. Hornung, "Enhancing scalability of parallel structured amr calculations," in *ICS '03: Proceedings of the 17th annual international conference on Supercomputing*. New York, NY, USA: ACM Press, 2003, pp. 336–347.
- [8] C. Burstedde, O. Gattas, G. Stadler, T. Tu, and L. C. Wilcox, "Towards adaptive mesh pde simulations on petascale computers," in *TeraGrid 08*, 2008.
- [9] C. Burstedde, O. Ghattas, M. Gurnis, G. Stadler, E. Tan, T. Tu, L. C. Wilcox, and S. Zhong, "Scalable adaptive mantle convection simulation on petascale supercomputers," in *SC '08: Proceedings of the 2008 ACM/IEEE conference on Supercomputing*. Piscataway, NJ, USA: IEEE Press, 2008, pp. 1–15.
- [10] H. Childs, E. S. Brugger, K. S. Bonnell, J. S. Meredith, M. Miller, B. J. Whitlock, and N. Max, "A contract-based system for large data visualization," in *Proceedings of IEEE Visualization 2005*, 2005, pp. 190–198.
- [11] J. Guilkey, J. Hoying, and J. Weiss, "Modeling of multicellular constructs with the material point method," *Journal of Biomechanics*, vol. 39, pp. 2074–2086, 2007.
- [12] A. Brydon, S. Bardenhagen, E. Miller, and G. Seidler, "Simulation of the densification of real open-celled foam microstructures." *J. Mech. Phys. Solids*, vol. 53, pp. 2638–2660, 2005.
- [13] G. Krishnamoorthy, S. Borodai, R. Rawat, J. Spinti, and P. Smith, "Numerical modeling of radiative heat transfer in pool fire simulations." Orlando, Florida: ASME International Mechanical Engineering Congress (IMECE), 2005.
- [14] B. Kashiwa, M. Lewis, and T. Wilson, "Fluid-structure interaction modeling," Los Alamos National Laboratory, Los Alamos, Tech. Rep. LA-13111-PR, 1996.
- [15] B. Kashiwa, "A multifield model and method for fluid-structure interaction dynamics," Los Alamos National Laboratory, Los Alamos, Tech. Rep. LA-UR-01-1136, 2001.
- [16] B. Kashiwa and E. Gaffney, "Design basis for cfdlib." Los Alamos National Laboratory, Los Alamos, Tech. Rep. LA-UR-03-1295, 2003.
- [17] D. Sulsky, Z. Chen, and H. Schreyer, "A particle method for history dependent materials," *Comput. Methods Appl. Mech. Engrg.*, vol. 118, pp. 179–196, 1994.
- [18] J. Guilkey, T. Harman, and B. Banerjee, "An eulerian-lagrangian approach for simulating explosions of energetic devices," *Computers and Structures*, vol. 85, pp. 660–674, 2007.
- [19] Y. Guo and J. Nairn, "Calculation of j-integral and stress intensity factors using the material point method." *Computer Modeling in Engineering and Sciences*, vol. 6, pp. 295–308, 2004.
- [20] D. Sulsky, S. Zhou, and H. Schreyer, "Application of a particle-in-cell method to solid mechanics," *Computer Physics Communications*, vol. 87, pp. 236–252, 1995.
- [21] J. Brackbill and H. Ruppel, "Flip: A low-dissipation, particle-in-cell method for fluid flows in two dimensions," *J. Comp. Phys.*, vol. 65, pp. 314–343, 1986.
- [22] —, "Flip: A method for adaptively zoned, particle-in-cell calculations of fluid flow in two dimensions." *Journal of Computational Physics*, vol. 65, pp. 314–343, 1986.
- [23] J. Brackbill, "Particle methods," *International Jour. Numer. Meths. in Fluids*, vol. 47, pp. 693–705, 2005.
- [24] J. Luitjens, "Malloc trace users guide." [Online]. Available: <http://www.csafe.utah.edu/wiki/index.php/Documentation/UsersGuide/MallocTrace>
- [25] G. Karypis and V. Kumar, *MeTis: Unstructured Graph Partitioning and Sparse Matrix Ordering System, Version 2.0*.
- [26] C. Walshaw, M. Cross, M. G. Everett, and S. Johnson, "Jostle: Partitioning of unstructured meshes for massively parallel machines," in *Parallel Computational Fluid Dynamics: New Algorithms and Applications*. Elsevier, 1994.
- [27] K. Devine, B. Hendrickson, E. Boman, M. S. John, and C. Vaughan, "Design of dynamic load-balancing tools for parallel applications," in *ICS '00: Proceedings of the 14th international conference on Supercomputing*. New York, NY, USA: ACM, 2000, pp. 110–118.
- [28] E. Boman, K. Devine, L. A. Fisk, R. Heaphy, B. Hendrickson, C. Vaughan, U. Catalyurek, D. Bozdog, W. Mitchell, and J. Teresco, *Zoltan 3.0: Parallel Partitioning, Load-balancing, and Data Management Services; User's Guide*, Sandia National Laboratories, Albuquerque, NM, 2007, tech. Report SAND2007-4748W.
- [29] J. Luitjens, M. Berzins, and T. Henderson, "Parallel space-filling curve generation through sorting: Research articles," *Concurr. Comput. : Pract. Exper.*, vol. 19, no. 10, pp. 1387–1402, 2007.
- [30] Q. Meng, J. Luitjens, and M. Berzins, "A comparison of load balancing algorithms for amr in uintah," University of Utah, SCI Technical Report UUSCI-2008-006, 2008.
- [31] D. C. Montgomery, L. A. Johnson, and J. S. Gardiner, *Forecasting and Time Series Analysis*, 2nd ed. Koga: McGraw-Hill, 1990.
- [32] R. E. Kalman, "A new approach to linear filtering and prediction problems," *Transactions of the ASME—Journal of Basic Engineering*, vol. 82, no. Series D, pp. 35–45, 1960.