

Assessing the Effects of Data Compression in Simulations Using Physically Motivated Metrics

Daniel Laney
Lawrence Livermore Lab
dlaney@llnl.gov

Steven Langer
Lawrence Livermore Lab
langer1@llnl.gov

Christopher Weber
Lawrence Livermore Lab
weber30@llnl.gov

Peter Lindstrom
Lawrence Livermore Lab
pl@llnl.gov

Al Wegener
Simplify
awegener@simplify.com

ABSTRACT

This paper examines whether lossy compression can be used effectively in physics simulations as a possible strategy to combat the expected data-movement bottleneck in future high performance computing architectures. We show that, for the codes and simulations we tested, compression levels of 3–5X can be applied without causing significant changes to important physical quantities.

Rather than applying signal processing error metrics, we utilize physics-based metrics appropriate for each code to assess the impact of compression. We evaluate three different simulation codes: a Lagrangian shock-hydrodynamics code, an Eulerian higher-order hydrodynamics turbulence modeling code, and an Eulerian coupled laser-plasma interaction code. We compress relevant quantities after each time-step to approximate the effects of tightly coupled compression and study the compression rates to estimate memory and disk-bandwidth reduction. We find that the error characteristics of compression algorithms must be carefully considered in the context of the underlying physics being modeled.

Categories and Subject Descriptors

E.4 [Coding and Information Theory]: Data compaction and compression; H.3.4 [Information Storage and Retrieval]: Systems and Software—*performance evaluations (efficiency and effectiveness)*; G.1.0 [Numerical Analysis]: General—*error analysis*

Keywords

data compression, high performance computing

1. INTRODUCTION

Data movement on future Exascale (10^{18} floating point operations per second) systems is expected to consume a substantial portion of their power budget [15]. At the same

time, a large number of applications will be limited by memory bandwidth rather than processor performance [21, 26]. These two factors mean that the “data movement bottleneck” on future systems will make it difficult to fully exploit Exascale HPC resources. This paper examines whether compression has the potential to be used effectively in physics simulations to mitigate these effects by reducing data size and therefore bandwidth requirements. We investigate the effects of compression on simulation results, a necessary first step before studying the kinds of algorithms and hardware implementations that would make compression feasible on future HPC resources.

In this paper we make a distinction between two possible uses of compression: “disk compression” and “memory compression”. In “disk compression,” arrays are compressed before they are written to disk; in “memory compression,” arrays are compressed while stored in off-chip memory and decompressed when they are loaded into on-chip cache for computation. We focus mainly on memory compression in this paper, approximating its effects by compressing and decompressing simulation state at the end of every cycle. Acceptable disk compression ratios will be at least as high as acceptable memory compression ratios.

1.1 Lossless vs. Lossy Compression

Compression provides the opportunity to trade off additional computation in exchange for reduced bandwidth. While lossless compression would enable this trade-off without modifying simulation results, lossless compression ratios are modest for the codes we consider and thus will only have modest impact on memory bandwidth requirements at Exascale. For example, we show in section 5.2 that data from LLNL’s pF3D code can only be losslessly compressed by a factor of 1.29X by *fpzip*, one of the floating point compression schemes used in this paper. Therefore, to have a significant impact on the memory or disk bandwidth requirements of simulations, lossy compression methods are necessary.

The main challenge when applying lossy compression is assessing its effect on simulation accuracy. Lossy compression is commonly used for photographs and videos where errors need only meet human perceptual requirements. Key-frames in videos periodically reset the error to zero. The physics codes we consider solve time-dependent differential equations where errors may accumulate over time. The changes in the solution due to the use of compression may grow more rapidly if compression errors at one time step are strongly correlated with errors at the next time step. The changes

due to the use of compression are thus likely to be functions of both the compression ratio and the error characteristics of the compression algorithm.

The fact that lossy compression may change simulation results is not necessarily a showstopper – *many decisions made in setting up a simulation also change the answer*. For example, scientists may choose a mesh resolution based on a trade-off between more accurate answers and the computational cost of the simulation. Computational scientists usually make these choices based on the desired accuracy of various integral physical quantities, not by comparing per-grid-point differences like the mean square error between two proposed simulations. We propose using the same integral physical quantities to assess the impact of compression.

1.2 Approach

In our approach, the code developer runs test problems with and without compression and evaluates key physical quantities. The results can be used to select the best compression algorithm and the proper compression level. We evaluate the impact of compression on three different simulation codes: LULESH, a shock hydrodynamics mini-app; Miranda, a hydrodynamics code for large turbulence simulations; and pF3D, which simulates the interaction of a high intensity laser and a plasma (ionized gas). These codes provide three different regimes in which to assess the effects of compression: LULESH is a single-physics code that models a stable problem; Miranda is also a single-physics code, but uses higher order methods and models unstable flows that can be sensitive to numerical errors; Finally, pF3D is a multi-physics code that couples two different physics models.

For each simulation code, we present a high-level discussion that should be accessible to a general HPC audience and discuss relevant physical metrics. We then present the effect of increasing compression ratios on these metrics. This paper shows that, in the simulation codes we examined, 3–5X compression ratios can be applied without causing a significant change in the physical quantities that are of most interest for the simulation.

2. PRIOR ART

Numerical compression has a rich history dating back to the 1960s. The compression research community has understandably focused on compression of consumer speech, audio, images, and video, given the ubiquity of such media. However, media compression techniques, which tend to be lossy, are inappropriate for a universal numerical encoder because the quality metrics for media compression are determined by limitations of human hearing and vision, not by the accuracy requirements of numerical computations.

Prior work on compressing scientific data has mostly focused on lossless compression [2, 8, 9, 20, 27], in part due to a poor understanding of the effects of lossy compression on the simulation results. However, lossless compression typically achieves less than a 2X reduction on floats and even lower compression ratios on doubles, greatly limiting its impact on data movement.

Lossy compression for scientific and medical data has historically been applied to data sets to reduce their size for post-processing and visualization. Muraki et al. [24] pioneered the application of wavelet transforms to volumetric data for visualization applications, thresholding wavelet coefficients to obtain approximate representations and apply-

ing non-orthogonal wavelets to medical images [25]. More recently, Woodring et al. [35] explored the use of wavelet compression in JPEG2000 to enable scientists to trade off data quality for reduced data size. Their approach targeted reduced disk I/O time in visualization applications and faster transfers of data from remote systems to local assets. As an alternative to transform coding, compressive sensing [4, 7] has the potential to rapidly reduce a stream of numerical data to any user-specified size, but requires significant back-end complexity to recover (an approximation of) the original signal with no direct control over accuracy.

ISABELA [16] has been proposed as a data reduction approach for HPC applications, particularly targeted at data deemed effectively incompressible due to randomness or noise. The data is partitioned into blocks and the data in each block is sorted by value. Splines are fitted to the resulting monotonic data using a relatively small number of knots. Temporal patterns are leveraged to further increase compression ratios. The approach has several advantages for disk compression, particularly since the method does not require global communication and can be run in-situ. Lehmann et al. [18] present a modified version of ISABELA suited to compressing flow simulations in porous media.

To the best of our knowledge, prior lossy compression techniques have not been considered within the context of memory compression, in which the simulation state itself is compressed. The simulations presented in this paper operate repeatedly on data that has been previously compressed, leading to the possibility of amplification of compression-induced differences. In the next section we discuss compression algorithms in this context, and motivate our choice to use predictive coders in our experiments.

3. COMPRESSION ALGORITHMS

In most physics simulations, the low order bits of floating point numbers are effectively random. The presence of these random bits prevents lossless methods from achieving high compression rates. If some of the low order bits of floating point numbers can be removed (or approximated) without a significant impact on physics answers, we can achieve greater compression at the expense of lost information. We believe a compression method suitable for in-memory compression should be able to run in both lossless and lossy modes, and to be amenable to implementation in hardware with minimal resources and power consumption.

We classify lossy compression schemes into three broad categories: general lossless schemes like *gzip* and *bzip2*, transform coding approaches such as those based on wavelet or discrete cosine bases, and predictive coding schemes like the two methods we apply in this paper. We do not consider general lossless compression methods in this paper as they are limited in the compression rates that can be achieved and tend to be less effective for floating point data. Although transform coding is traditionally used for lossy compression, these methods have two main drawbacks: First, bit-for-bit lossless compression with these schemes can be difficult to achieve due to the subtleties of floating-point arithmetic, including rounding modes and error, catastrophic cancellation, order-of-evaluation dependence, extended precision, etc. Second, hardware implementations for the in-memory use-case are likely to use too many resources to be practical.

For these reasons, and for speed, we focus on predictive coding approaches, in which each floating-point number is

predicted based on a trend of recently encoded values. The prediction *residual* (difference between actual and predicted value) tends to be small and can generally be encoded using fewer bits than the original floating-point value. The number of recent values is generally quite small for a 1D predictor and therefore a hardware implementation can be quite compact. Typically, a lossless compression scheme is then used to compress the residuals so that the overall approach is bit-for-bit lossless. Lossy compression can be achieved by ignoring low order bits of the input values, resulting in smaller residuals and higher compression rates. For example, applying lossless compression after zeroing the bottom 32 bits of a double precision number usually produces total compression ratios around 3–4X. Finally, the errors introduced by predictive schemes are independent of neighboring values, usually being the result of truncation or rounding which occur already in fixed precision floating point computations, whereas the errors introduced in transform coding (e.g., Gibbs ringing) have spatial extent and exhibit spatial correlations.

In this paper we study two compression algorithms designed for floating-point data: Simplify’s *APAX* (APplication AXceleration) encoder [30, 31] and the *fpzip* [19, 20] compressor developed at LLNL. The *APAX* algorithm has previously been applied to climate data [10], computed tomography x-ray samples [33], and a variety of integer and floating-point datasets [32]. We describe these two compressors in parallel to highlight their similarities and differences.

3.1 The *fpzip* and *APAX* Compressors

When used for lossy compression, *APAX* and *fpzip* both begin by quantizing the value f to be encoded, in effect reducing its precision. Quantization in *APAX* involves converting a *block* of N consecutive values (usually $N = 256$) to a signed integer representation, which can be thought of as aligning the floating-point values in a block to a common largest exponent. If exponent differences are large, this uniform quantization step may result in some loss of precision for the smallest (in magnitude) values in the block. After exponent alignment, each signed value is treated as a 32-bit integer \hat{f} . For (64-bit) double-precision data, this implies that the bottom 32 bits of the significand are discarded. This design decision was motivated by performance reasons – by modifying *APAX* to use 64-bit integer arithmetic, such truncation could be avoided, possibly at the expense of slower compression. *APAX* then quantizes the integer \hat{f} uniformly to $\bar{f} = \text{round}(\hat{f}/q)$, where q is a quantization level either specified by the user or computed adaptively by *APAX* to meet a target coding rate.

fpzip delays the integer conversion and leaves the values in their floating-point representation, quantizing the significand instead. *fpzip* restricts q to be an integer power of two, which effectively leads to truncation of the significand by discarding (zeroing) some fixed number of least significant bits. Setting $q = 1$ guarantees entirely lossless compression. This non-uniform quantization allows the *relative error* to be bounded in *fpzip*, whereas within each *APAX* block the *absolute error* is bounded. Quantization is the only potential source of loss in both compressors.

Following quantization, each value is predicted as a linear combination of recently encoded values. Both compressors rely on polynomial interpolation for prediction, with fixed integer polynomial coefficients. *APAX* uses univariate La-

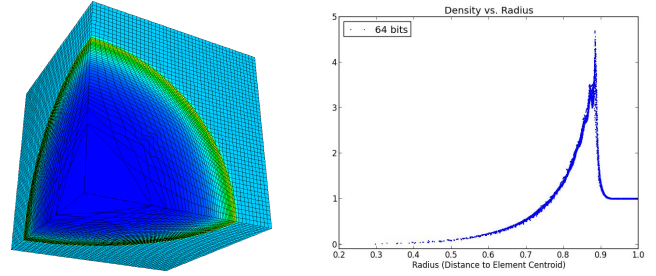


Figure 1: An overview of the LULESH shock-hydrodynamics simulation. (left) Density field showing the shock wave and mesh deformation. (right) Scatter plot of density vs. radial displacement from the origin for every element in the mesh.

grange polynomials of degree 0 and 1, allowing linear polynomials (or any function with $\frac{\partial^2 \bar{f}}{\partial x^2} = 0$) to be reproduced. That is, the unknown value $\bar{f}[x]$ is predicted in terms of the n previous, known values by solving $\sum_{i=0}^n (-1)^i \binom{n}{i} \bar{f}[x-i] = 0$. The “best” polynomial degree $n-1$ is chosen locally by monitoring its effect on compression.

fpzip exploits correlations in more than one dimension using the Lorenzo predictor [11], which over a 3D domain reproduces trivariate quadratic polynomials (or any function for which $\frac{\partial^3 f}{\partial x \partial y \partial z} = 0$). In 3D this predictor solves $\sum_{i,j,k \in \{0,1\}} (-1)^{i+j+k} f[x-i, y-j, z-k] = 0$ for $f[x, y, z]$, and thus requires buffering a whole 2D “slice” from the 3D domain. Note that this predictor uses only additions and subtractions of known values.

Given the true floating-point value f and its prediction p , *fpzip* converts f and p to integers \bar{f} and \bar{p} via a monotonic mapping that treats the binary floating-point number as a sign-magnitude integer. This step is not needed in *APAX*, where the integer conversion occurs earlier. Both compressors then compute an integer residual $\bar{r} = \bar{f} - \bar{p} = s(2^e + d)$ with sign $s \in \{-1, 0, +1\}$, exponent e , and e -bit difference $d \in \{0, \dots, 2^e - 1\}$. The bits of d generally exhibit no correlation and are transmitted verbatim. However, the sequence of “signed exponents” $\tilde{e} = s(e+1)$ tends to be highly correlated (if not necessarily peaked around zero). *APAX* exploits spatial correlations by encoding differences between consecutive exponents in small groups. *fpzip*, on the other hand, models the non-uniform distribution of exponents and encodes each \tilde{e} independently using a fast entropy coder.

4. SIMULATION & VALIDATION METHODS

4.1 LULESH

LULESH is a shock hydrodynamics mini-application developed for use in evaluation of current and future computer systems and proposed programming models [12–14]. LULESH solves the Sedov blast wave test problem – a point explosion surrounded by an initially uniform surrounding gas. The gas consists of one material and is modeled in three dimensions using a Lagrangian (moving mesh) formulation. The point energy source for the explosion is approximated by initializing the zone at the origin with a large initial energy. Figure 1 depicts two key features of the simulation – the deformation of mesh elements and the shock wave.

LULESH solves the inviscid compressible Navier-Stokes equations in the Lagrangian formulation. A staggered mesh approximation [34] with single point quadrature for element-centered thermodynamic quantities such as density and pressure is utilized. Kinematic variables such as position and velocity are defined at mesh nodes. The Sedov problem presents an interesting use case for compression, as most field values vary slowly over most of the domain, but change quickly and abruptly near the shock.

For the Sedov problem, the two key physics requirements are that the blast wave should be spherical, and the shock radius versus time should match the analytic solution. In this paper we evaluate the symmetry of the shock by comparing the field values as a function of radius for various compression levels and methods, to the results obtained with a double precision simulation with no compression. We assess physical accuracy by estimating shock position as the distance to the centroid of the element with maximal density, and comparing that value between runs with and without compression. Although this measure is not the most sophisticated way to assess shock position, it is simple and can be consistently applied across LULESH runs at varying mesh resolutions. Finally, we measured differences in internal energy between compression and non-compression runs to validate that compression wasn't violating the expected behavior of the LULESH simulation as represented by full double precision runs at different mesh resolutions.

4.2 pF3D

The National Ignition Facility (hereafter NIF; [22, 23]) is an NNSA experimental facility that houses the world's most powerful laser. One of the key goals of the NIF is to compress a target filled with deuterium and tritium to a temperature and density high enough that fusion ignition occurs.

The intensity of NIF beams exceeds 10^{15} W/cm² in the brightest spots. When intensities are this high, it is possible for the laser to couple to fluctuations in the plasma density and backscatter a significant fraction of the laser light. pF3D [1, 17, 28] is a multi-physics code that simulates interactions between laser beams and the plasma in NIF experiments. pF3D is used to evaluate proposed target designs and to pick the ones with acceptably low levels of backscatter.

pF3D zones are roughly the size of the laser wavelength ($0.35\mu\text{m}$) while the plasmas of interest are several mm across. Simulations of the full path of a single laser beam require 50 billion or more zones. Simulations of five interacting beams may require more than a trillion zones. pF3D has been used to run simulations with 64k or more cores on IBM Blue Gene/L, IBM Blue Gene/P, and Cray XE-6 systems. Work is underway to run a million core simulation on an IBM Blue Gene/Q system.

Figure 2 shows the intensity pattern of the laser beam as it enters the target in a NIF experiment. The bright spots are referred to as speckles. Speckles are narrow transverse to the laser propagation direction and extend for many laser wavelengths in the propagation direction. The laser is designed so that the speckles in the beam move around in time. The plasma temperature and density respond to the intensity averaged over time, so they see a smooth beam.

Backscattered light can be generated in the brightest speckles and grows in strength faster than the speckles move. The laser light will refract when there is a density gradient. Speckles are long enough that a small change in density can

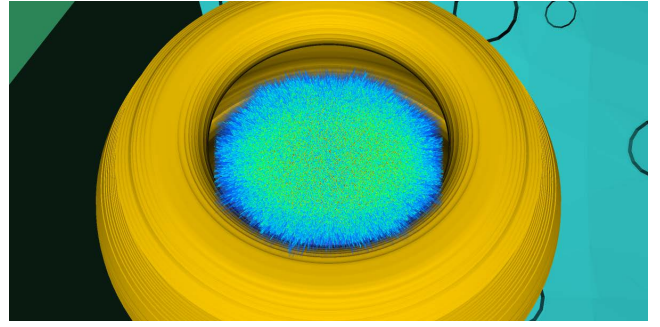


Figure 2: The figure shows the intensity of the laser beam in a pF3D simulation as it enters the hohlraum. The hohlraum is a can-like object with the fusion target at its center. The intensity increases as the color changes from blue to green to red. The bright spots embedded in a lower intensity background are a design feature of NIF beams.

shift a speckle sideways by a significant fraction of its width.

The locations of the speckles are extremely sensitive. For example, the fused multiply-add instruction of the PowerPC processor (with internal registers that are longer than 64 bits) produces results sufficiently different from x86_64 processors to cause speckles to move sideways by a zone or more. Compiler optimizations that change the order in which summations are carried out also cause speckles to move.

As a result of this high sensitivity to such small variations, validation of new versions of pF3D relies on comparisons of the total transmitted and reflected light as a function of time (the total light is insensitive to the exact placement of speckles) and on intensity histograms. This paper uses the same physics-based validation methods to assess the impact of lossy compression on pF3D.

4.3 Miranda

Miranda [3, 6] is a Navier-Stokes code used to simulate a range of hydrodynamic problems with higher-order accuracy. It uses spectral methods or compact differencing to resolve turbulent structures with minimal dissipation. Dissipation is added at high wavenumbers through the use of artificial fluid properties [5], which act as a large-eddy simulation (LES) subgrid model. Miranda has run simulations using over 64k cores on Blue Gene/L systems.

The Miranda test problem simulates the growth of the Rayleigh-Taylor instability (RTI). Simulations of the RTI typically start with small-scale perturbations on the interface separating two fluids of different densities. The high density fluid is on top of the low density fluid – an unstable situation. The perturbation amplitudes grow, neighboring perturbations merge, and, eventually, turbulent mixing occurs. This inverse cascade of scales, from the initial short wavelength perturbation to large wavelengths at late times, requires high-order accuracy to ensure relevant features are not removed through dissipative numerics or influenced by the amplification of numerical noise.

The test problem is physically unstable to all perturbation wavelengths. Filtering is employed to damp the growth of modes with short wavelengths. The damping is quite strong at a wavelength of 2 zones, but drops quickly as the wavelength increases. An important point to remember when

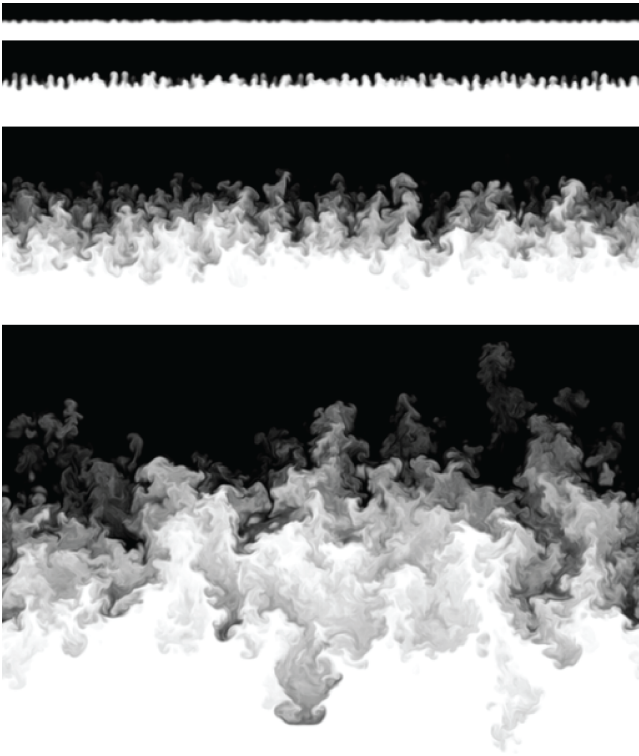


Figure 3: Density fields are shown from a Miranda simulation of the Rayleigh-Taylor instability at (from top to bottom) $t/\tau = 1, 2.5, 10, 20$.

investigating compression schemes is that the physics will amplify numerical perturbations introduced by compression as well as deliberately imposed perturbations. This means that the wavelength spectrum and step-to-step coherence of the perturbations produced by compression matter, in addition to the compression ratio.

The incompressible simulations shown here are initialized with narrowband Gaussian perturbation spectra peaked at 8 grid cells per wavelength (λ_0) and an RMS amplitude of 0.1 grid cells. This allows the instability to begin in its linear stage. The two fluids have densities of $\rho_1 = 1$ and $\rho_2 = 3$, and gravitational acceleration of $g = 1$, providing a time scaling of $\tau = \sqrt{\lambda/Ag}$, where $A = (\rho_2 - \rho_1)/(\rho_2 + \rho_1) = 0.5$ is the Atwood number. Density fields from a $512^2 \times 1024$ run using 1024 processors are shown in Figure 3. In the top two images, at $t/\tau = 1$ and 2.5, the perturbations – just barely noticeable in the first image – are growing independently at an exponential rate. By $t/\tau = 10$, shown in the third image, perturbations have merged, producing larger scales and mixed fluid. At $t/\tau = 20$, shown in the lower image, the layer has entered into an apparent turbulent state.

An important quantity in assessing RTI simulations is the mixing layer thickness, h , which is expected to behave as $h^2 = 4\alpha Agh$ when the layer reaches a self-similar state at late times. Another important characteristic at late times is the spectrum of perturbations as a function of spatial frequency. These physically motivated quantities will be used to assess the differences between compressed and uncompressed simulations.

5. RESULTS

This section presents results on the acceptable level of compression in the LULESH, Miranda, and pF3D test problems. Simulations with and without compression are compared using the physics metrics mentioned in the earlier discussion. Simulations may be sensitive to the details of how compression is performed. As a result, the acceptable level of compression for a given test problem may differ between different compression algorithms.

The codes have been modified so that they compress and decompress variables at the end of each time step. This procedure approximates “memory compression,” which, if actually in place between the memory and the cache, would occur several times per time step. We also ran a Miranda test where the compression function was called at the lower frequency of checkpoint dumps to approximate “disk compression”.

We report the minimum, maximum, and average compression ratios over all domains. The minimum compression ratio is the most important measure when using “memory compression.” The process with the lowest compression ratio will have the hardest time fitting in the available memory, and it will spend the most time reading and writing memory. All three of these codes use a bulk synchronous programming model, so the slowest process controls the performance. Efficient parallel I/O packages have performance that is dependent on the total number of bytes written, not on the number contributed by each process. With an I/O package of this sort, the average compression ratio is the important quantity.

5.1 LULESH

We explore the effects of compression on the accuracy of the LULESH Sedov blast wave simulation by investigating the interaction between compression and mesh resolution, and whether accumulation of errors occurs in longer running simulations. We performed two studies: in the first study we varied the mesh resolution and ran *fpzip* over a full range of precisions, to $t = 0.01$. We kept the total energy in the problem constant by scaling the initial energy density of the zone at the origin with mesh resolution. In the second study, we ran with both *APAX* and *fpzip* on a 105^3 mesh to time $t = 0.05$ at a small set of target compression rates (*APAX*) or precisions (*fpzip*). In the second study the energy at $t = 0$ was adjusted so that the shock wave traversed nearly to the boundary. We measured the shock position in a manner similar to Tasker et al. [29]. Figure 4(left) shows density vs. radius for all elements in a 105^3 simulation, with the shock position taken as the radius at the maximum density value in the plot.

In the first study, we ran a series of simulations with *fpzip* at precisions from 24 to 64 bits in steps of 4 bits on several mesh resolutions. Figure 4(middle) shows that the shock position is stable down to 28 bits of precision, and that the differences are smaller than those caused by changing the mesh resolution of this unconverged problem. Figure 4(right) shows that the root mean square (hereafter RMS) error varies with compression rate and does not exhibit the abrupt transition between 28 and 24 bits seen in the middle plot. Compression results were largely independent of mesh resolution in this study. Our assessment is that with respect to shock position the LULESH Sedov simulation has correct behavior down to 48 bits, and that in several

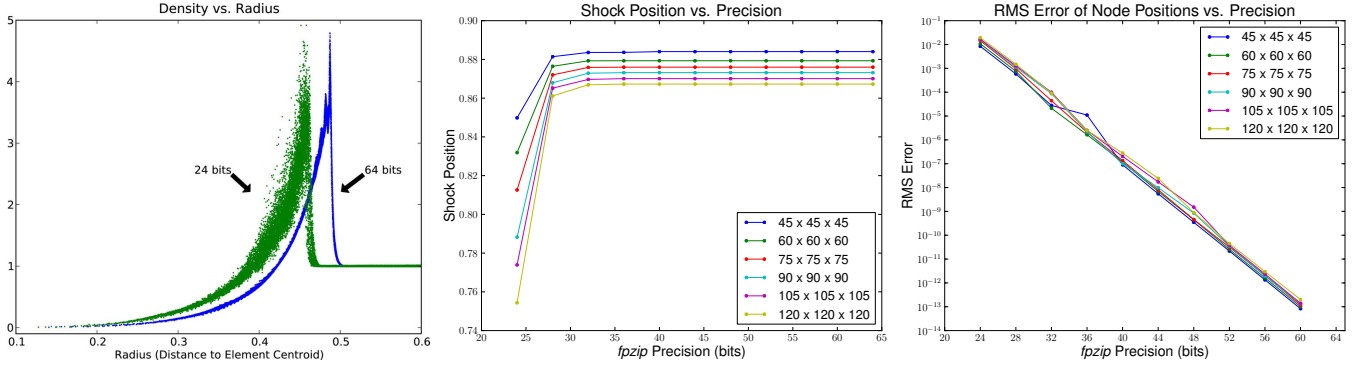


Figure 4: (left) Density plotted against radius at $t = 0.02$, for *fpzip* precision values of 24 and 64 bits, and a 105^3 mesh size. When retaining only 24 bits of precision, a significant amount of energy is lost, causing the shock to be delayed. (middle) Shock position as a function of mesh size and *fpzip* precision. Variations in shock position between 28 and 64 bit *fpzip* are smaller than the variation at fixed compression between 45^3 and 120^3 mesh resolutions. (right) The root mean square error of nodal positions varies smoothly with *fpzip* precision rather than exhibiting the abrupt transition seen in the middle plot. The dependence on mesh resolution is weak.

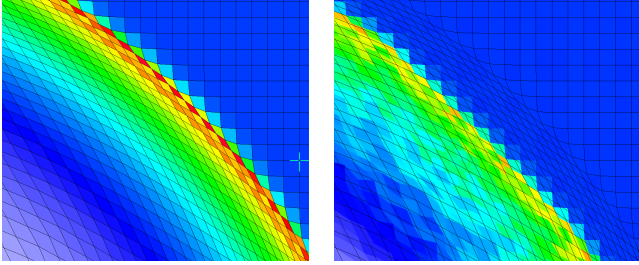


Figure 5: (left) A zoom in of the density field at the shock in a 105^3 mesh at full precision. (right) The shock at 24 bits of precision, showing that a significant amount of noise is introduced, although this noise is symmetric. In both images, blue indicates low density, and red indicates the highest density.

aspects precisions as low as 32 bits are sufficient. We confirmed this assessment by measuring internal energy loss in longer-running simulations as we present below.

In the second study we ran simulations at 105^3 resolution to $t = 0.05$ (4304 time steps), in order to assess whether there was error accumulation over time due to many applications of compression. This study showed that the shock position was accurately captured at *fpzip* precisions down to 32 bits, and *APAX* rates up to 4X. The density vs. radius plot in Figure 4(left) shows density vs. radius for *fpzip* precisions of 24 and 64 bits (lossless), showing a failure mode when the precision is too low. The shock radius is smaller than in the higher precision runs, indicating that energy has been lost and the shock is trailing the full-precision solution. Figure 5 shows zoomed in views of the density field at the two precisions showing that the shock is no longer spherical. We note that *fpzip* may introduce a bias as it first discards significant bits, effectively rounding towards zero. However, we found that at high compression rates *APAX* also loses energy, presumably due to the scaling process it uses to convert floating point values into 32-bit integers before compressing them.

compressor	ratio	energy	result
<i>fpzip</i> 64-bit	2.0 / 2.5 / 7.3	0.0	pass
<i>fpzip</i> 48-bit	2.8 / 3.7 / 9.7	$-1 \times 10^{-6}\%$	pass
<i>fpzip</i> 32-bit	3.8 / 6.5 / 14.7	-0.07%	pass
<i>fpzip</i> 24-bit	N/A	N/A	fail
<i>APAX</i> 2X	2.4 / 3.1 / 19.2	$+6 \times 10^{-6}\%$	pass
<i>APAX</i> 3X	3.0 / 3.5 / 19.2	-0.002%	pass
<i>APAX</i> 4X	4.0 / 4.5 / 19.0	-0.01%	see text

Table 1: This table shows the compression ratios achieved for the LULESH runs for mesh size 105^3 . The ‘ratio’ column shows the minimum, mean, and maximum compression ratios for all fields in the problem. The value for mean compression is the ratio between the storage required for all fields, divided by the total compressed size of all fields. The ‘energy’ column records the percent of internal energy lost (or gained), with respect to the full precision run. The compression algorithm and setting for each run is shown and is labeled by whether it passes the physics criteria.

Table 1 shows the compression ratios and final internal energy change achieved for LULESH runs using both *fpzip* and *APAX* compression for the second study. We see that over a fairly large range of compression rates and precisions that minimal energy is lost. The 4X *APAX* run showed good agreement with internal energy and shock position but seeded non-physical noise in the nodal positions ahead of the shock, so we did not deem that a total ‘pass’, nor a total ‘fail’. The Sedov problem is stable, and it appears that this noise in nodal positions in front of the shock does not impact the shock position adversely, at least up to the time we ended the simulation.

A limitation of our results is that the Sedov problem is a simple test case. The requirements of simulations involving more complicated geometries and multiple materials may have more restrictive requirements with respect to lossy compression. In addition, LULESH is a mini-app and contains fields whose sole purpose is to reproduce data move-

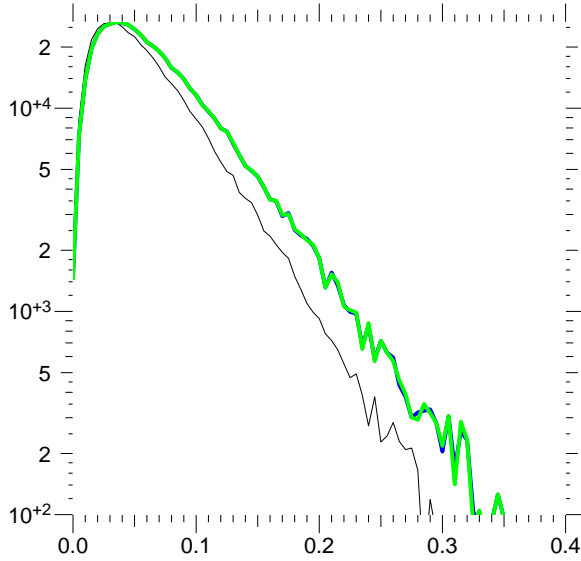


Figure 6: The figure shows a histogram of the amount of laser energy per bin as a function of the laser intensity. The beam has crossed the full extent of the plasma. The x-axis is scaled so that 1.0 is the peak intensity of the laser on the entrance plane. The black curve is from a run using 24-bit *fpzip*, the blue curve is from a run using 32-bit *fpzip*, and the green curve is from a run using uncompressed double precision variables. The blue and green curves overlay one another. The curves for runs using APAX compression are very similar and are not shown to avoid visual clutter.

ment patterns in larger codes. Since these fields are unit-valued and do not impact simulation results we do not include them in our compression ratio data. Finally, LULESH does not converge to a shock position cleanly, and up to mesh resolutions of 175^3 the shock position continued to vary slightly while keeping total energy constant.

5.2 pF3D

The pF3D test problem simulates the propagation of a laser beam through a plasma. The plasma is divided into $1 \times 2 \times 36$ equal domains – one for each of the 72 MPI processes. There are a total of $256 \times 1024 \times 756$ zones, each of which is $2 \times 2 \times 2$ wavelengths in size. Some of the laser light is absorbed, some is backscattered by Stimulated Brillouin Scattering (SBS), and some reaches the exit plane. When compression levels are high, some of the laser energy is “eaten by numerical gremlins.” The total amount of laser energy reaching the exit plane and the total amount of SBS crossing the entrance plane are key integral measures of what happens in a simulation. A more detailed metric is the amount of laser energy as a function of laser intensity.

Figure 6 shows a histogram of the amount of laser energy per bin as a function of the laser intensity. The histogram is made after the laser has crossed the full extent of the plasma, so energy has been lost due to absorption, backscattered light, and numerical effects. The x-axis is scaled so that 1.0 is the peak intensity of the laser on the entrance plane. All of the pF3D simulations used double precision arithmetic. The blue (32-bit *fpzip*) and green (lossless *fpzip*) curves are

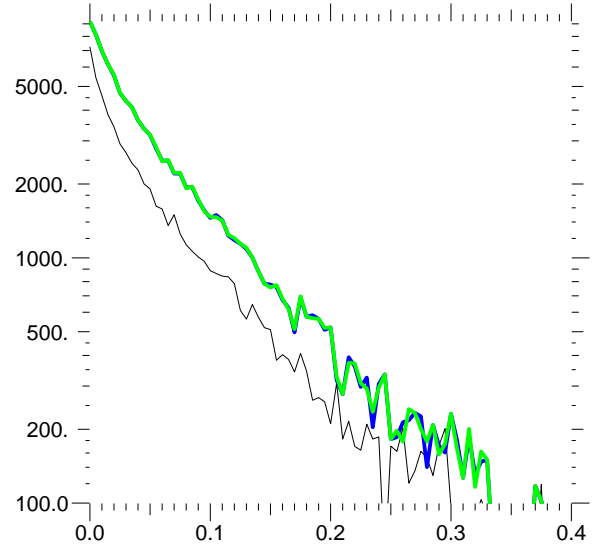


Figure 7: The figure shows a histogram of the amount of backscattered energy per bin as a function of the laser intensity. The histogram is made at the entrance plane where the backscattered light is at its maximum. The x-axis is scaled so that 1.0 is the peak intensity of the laser on the entrance plane. The black curve is from a 24-bit *fpzip* run, the blue curve is from a 32-bit *fpzip* run, and the green curve is from a run without compression. The blue and green curves overlay one another.

compressor	ratio	energy	result
<i>fpzip</i> 64-bit	1.27 / 1.29 / 1.32	0	pass
<i>fpzip</i> 48-bit	1.86 / 1.91 / 1.97	0	pass
<i>fpzip</i> 32-bit	3.46 / 3.66 / 3.89	-0.2%	pass
<i>fpzip</i> 24-bit	5.89 / 6.57 / 7.37	-37%	fail
APAX 2X	2.19 / 2.22 / 2.26	0	pass
APAX 3X	3.03 / 3.03 / 3.09	0	pass
APAX 4X	4.06 / 4.06 / 4.14	0.1%	pass
APAX 5X	5.06 / 5.07 / 5.16	0.9%	pass
APAX 6X	6.09 / 6.09 / 6.20	1.2%	pass
APAX 7X	7.10 / 7.11 / 7.22	10.7%	fail

Table 2: This table shows the compression ratios achieved for the pF3D runs. The ‘ratio’ column contains minimum, mean, and maximum compression rates. Each run is characterized as passing or failing the physics criteria, or crashing. The ‘energy’ column shows the percent difference in the backscattered energy relative to the uncompressed run.

effectively identical. The black curve (24-bit *fpzip*) is significantly lower at all but the lowest intensities and indicates that compression has led to the loss of laser energy. Compression is applied to the electric field of the laser, but the laser energy depends on the square of the field. A compression scheme that preserves the integral of the electric field can still significantly alter the laser energy. The compression scheme may also alter the direction in which portions of the laser beam travel. Laser energy propagating at an angle to the z -axis will travel a greater distance and suffer more absorption before reaching the exit plane.

Figure 7 shows a histogram of the amount of backscattered energy per bin as a function of the intensity at the entrance plane. The x -axis is scaled so that 1.0 is the peak intensity of the laser on the entrance plane. The blue (32-bit *fpzip*) and green (uncompressed) curves are effectively identical. The black curve (24-bit *fpzip*) has significantly lower energy at all intensities, a natural consequence of the reduced laser energy reaching the back of the simulation. APAX curves for compression levels that pass are very similar to the green and blue curves.

Table 2 shows the compression ratios achieved for the pF3D runs. The 64-bit *fpzip* run used lossless compression and only achieved a 23% reduction in data size. The 48-bit and 32-bit *fpzip* runs passed the physics criteria while the 24-bit *fpzip* run failed due to too much loss of energy in the transmitted laser light. The simulations using APAX compression passed the physics criteria up to 6X compression but failed at 7X due to too much increase in the transmitted laser light. We also ran pF3D with float (32-bit) precision and no compression. This run passed the physics criteria and could be thought of as achieving 2X compression. The run with 48-bit *fpzip* also achieved roughly 2X compression, but had higher precision than the run using float precision.

5.3 Miranda

To compare the effects of lossy compression on the development of the RTI, simulations with various levels of compression were run on $512^2 \times 1024$ grids using 1024 processors. After each time step, the density and the three components of velocity are compressed and decompressed using *fpzip* or APAX. This results in $\sim 10,000$ lossy compression steps over the course of the simulation.

In the case of *fpzip*, the compression step truncates the 64-bit data to either 48, 40, or 32 bits. Larger compression rates are achieved ($\sim 4X$) in regions where the flow field is relatively uniform (lossless compression is applied after the numbers are truncated). To test the effect of compressed checkpoint files, an additional test was run calling the 32-bit compression step every 500th time step, resulting in 16 compress/decompress calls over the course of the simulation. In contrast, APAX was used in three simulations with 2X, 4X, and 5X compression at each time step.

Images of the mixing layer at the end of the simulation ($t/\tau = 22$) are shown in Figure 8 from the 64-bit case (left) and the APAX 5X case (right). While there are differences between the two images, the quantitative metrics presented below show that the integral quantities and the turbulence state are nearly identical.

The thickness of the mixing layer, h , is plotted in Figure 9 as the lower group of lines. The reference simulation (64-bit) is plotted in black while the 40- and 32-bit *fpzip* compression cases (called each time step) are plotted in blue and red, respectively. The APAX 5X case is plotted as a green dashed line. The other four cases considered are all plotted in gray, since their differences with each other and with the 64-bit curve are minor. The simulations were stopped once the dominant wavelength approached the size of the simulation domain. The 32-bit *fpzip* case is the only one that failed outright, as it crashed at $t/\tau \approx 14$ when it became numerically unstable. The other cases, and even the 32-bit *fpzip* case before crashing, differed little from the reference simulation. At the time the 32-bit *fpzip* case crashed, its mixing-layer thickness was 2.4% larger than the reference

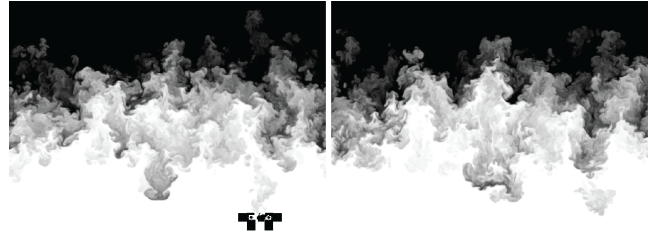


Figure 8: Density field at the end of the simulation from the reference (64-bit) run (left) and the APAX 5X case (right), which have nearly the same mixing-layer thickness and turbulence characteristics.

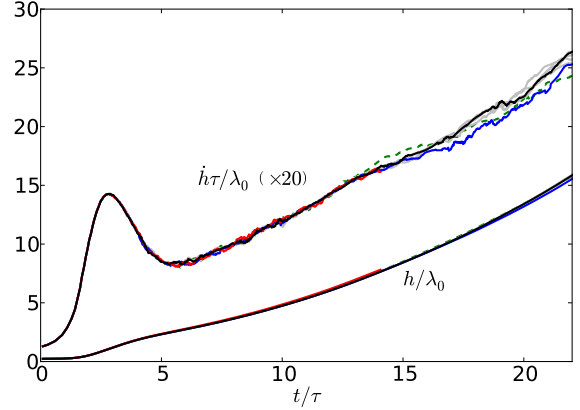


Figure 9: Mixing layer thickness (lower curves) and growth rate of the mixing layer (upper curves) from a reference (64-bit) calculation (black), an *fpzip* 32-bit compression case (red), an *fpzip* 40-bit compression case (blue), and an APAX 5X case (green dashed). The grey curves include a 48-bit *fpzip* case, a 32-bit *fpzip* case where compression is applied every 500th step, and APAX with 2X and 4X compression at every time step. Small differences are observable in the 32- and 40-bit *fpzip* cases. The 32-bit *fpzip* case crashed at $t/\tau = 14$.

simulation. The 32-bit *fpzip* case where compression was only applied every 500th time step ran without issue and differed by 0.38% in mixing-layer thickness at the end of the simulation ($t/\tau = 22$).

The growth rates are also shown in Figure 9 as the upper set of lines with the same color scheme. At early times, when the layer is growing exponentially, there are no differences between the seven curves. Once nonlinear growth begins, after $t/\tau > 3$, small differences are noticeable. These differences are minor and all cases show the expected $\dot{h} \propto t$ self-similar growth beyond $t/\tau > 6$. Near the end of the simulation the *fpzip* 40-bit case and the APAX 5X case had 5–7% smaller growth rates, resulting in 1.9% and 0.49% smaller mixing-layer thickness, respectively. Since these growth rates are systematically smaller, rather than simply having different random fluctuations like the other cases (gray curves), it is likely that the mixing-layer thickness difference will compound and become unacceptable if the simulation were run later in time (requiring a larger initial domain). The differences in mixing-layer thickness at the end of the simulation

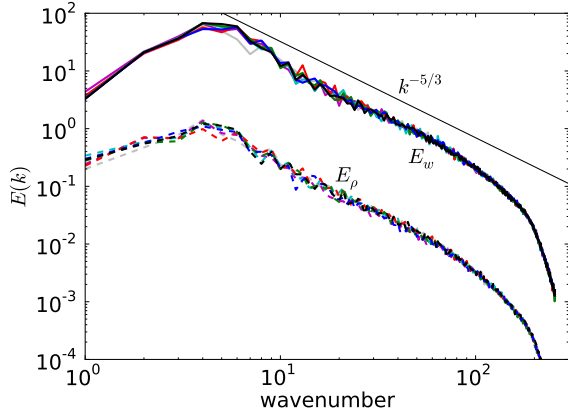


Figure 10: Energy spectrum of vertical velocity (solid) and density (dashed) at the mid-plane of the mixing layer at $t/\tau = 22$. All cases are shown except the 32-bit *fpzip* case, which did not reach this time. The spectra are very similar and feature a $k^{-5/3}$ inertial range.

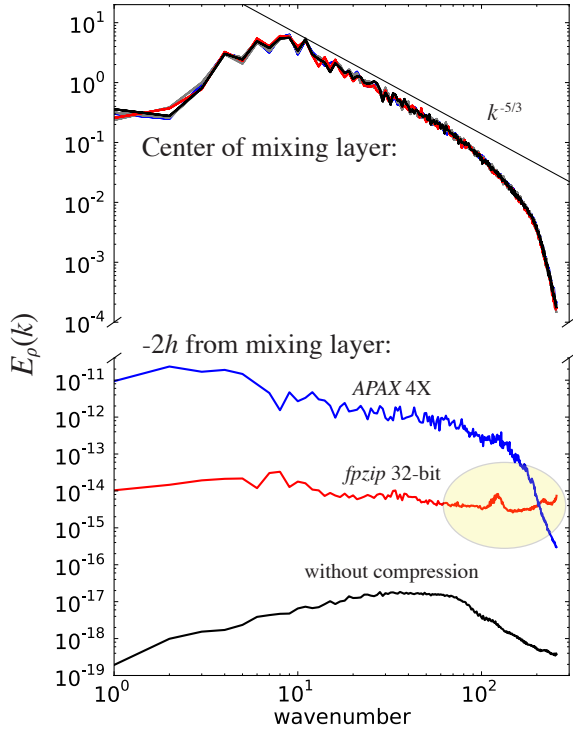


Figure 11: Energy spectrum of density at the mid-plane of the mixing layer (top set of curves) and at a distance $-2h$ from the mixing layer (bottom set of curves) at $t/\tau = 14$. The 32-bit *fpzip* case exhibits high wavenumber features that cause the problem to become numerically unstable.

are listed in Table 3.

The two-dimensional energy spectrum from the plane centered within the mixing layer is shown in Figure 10 at $t/\tau = 22$ for both the vertical velocity (solid) and the density (dashed) (the 32-bit case is not shown, as it did not reach

compressor	ratio	loss	thick.	result
<i>fpzip</i> 48-bit	1.8 / 3.3 / 4.6	5.9E-12	0.07%	pass
<i>fpzip</i> 40-bit	2.3 / 4.3 / 6.3	1.7E-09	1.9%	pass
<i>fpzip</i> 32-bit	3.2 / 4.5 / 6.5	4.1E-07	N/A	crash
APAX 2X	2.1 / 2.3 / 2.6	3.3E-10	0.007%	pass
APAX 4X	4.0 / 4.1 / 4.3	3.6E-06	0.19%	pass
APAX 5X	5.1 / 5.3 / 5.4	4.7E-05	0.49%	pass

Table 3: This table shows the compression ratios achieved for selected Miranda runs. The ‘ratio’ column shows the minimum, mean, and maximum compression rate over the simulation domain. The ‘loss’ column reports the time average of the RMS difference due to lossy compression divided by the mean density. The ‘thick.’ column shows the percent difference in mixing-layer thickness at the end of the simulation. Each run is characterized as passing or failing the physics criteria, or crashing.

this time). These spectra show the amount of energy contained in the various length scales within the problem. Most of the energy is in the large wavelengths (low wavenumbers). At the smallest scales (high wavenumbers) the energy is removed through viscous dissipation. In-between these scales is a power-law region, which, if the flow is turbulent, has a $-5/3$ slope. Aside from small fluctuations, the curves are nearly identical and all exhibit a $-5/3$ power-law for over a decade in wavenumbers, signifying turbulent behavior.

To further investigate why the *fpzip* 32-bit case crashed, Figure 11 shows density spectra at $t/\tau = 14$. The top set of curves show the energy spectra of the density field at the center of the mixing layer from all of the cases considered. At this location the spectra are all nearly identical, showing the beginnings of a power-law range and a dissipation region at high wavenumbers caused by Miranda’s LES filtering. The lower set of curves show the spectra from three simulations, taken from a plane located a distance $-2h$ from the center of the mixing layer, where the flow field is relatively quiescent. The APAX 4X case produced a greater loss of precision than the *fpzip* 32-bit case, which is reflected in its larger energy level in this region. In the *fpzip* 32-bit case, however, a bump can be seen in the high wavenumber region that corresponds to approximately 5 grid cells per wavelength. This scale is large enough to remain after Miranda’s filtering routines. This spectrum also shows an up-turn at the highest wavenumber portion of the spectrum. It is these features that accumulate and interact with Miranda’s numerics that cause the problem to crash. This shows that precision is not the only factor to consider when choosing a compression scheme for “memory compression” – one must also consider, or test, how compression will interact with the code’s numerics.

The compression rate varied over the course of the simulation, within different regions of the simulation domain, and for the different fields being compressed. In the APAX 5X case, for example, the three velocity fields were all compressed to 5–5.5X, while the density field started out at 14X compression in regions far from the mixing layer, and then reduced to $\sim 5X$ compression later in time. The rates were averaged across time and fields and are reported in Table 3. Also reported in this table is the loss that compression introduced at each time step. This is defined as the RMS

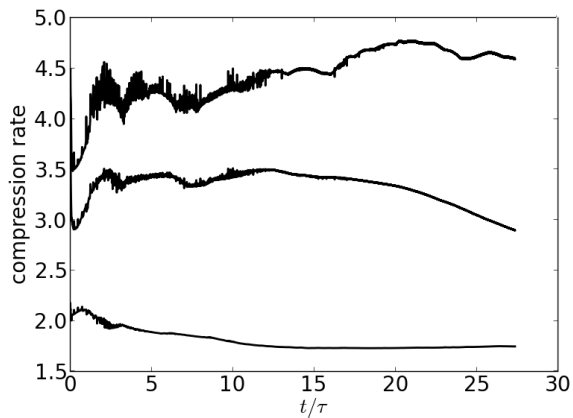


Figure 12: The compression ratio for an *fpzip* 48-bit compression run is shown. The three curves represent the maximum, minimum, and mean compression rate among the 1024 processes.

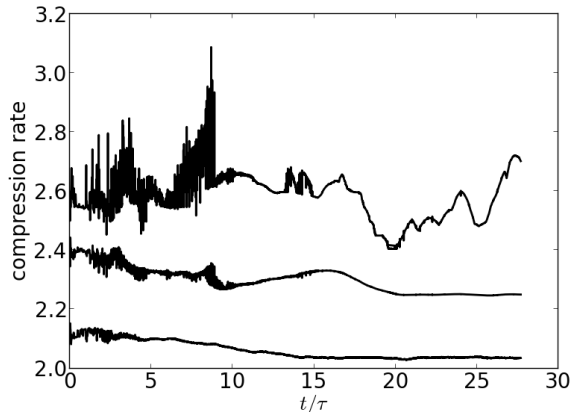


Figure 13: The compression ratio for a 2X *APAX* compression run is shown. The three curves represent the maximum, minimum, and mean compression rate among the 1024 processes.

difference in the density field after compression was applied and normalized by the mean density. This loss criterion was evaluated in the center of the mixing layer. As noted in the table, all of the runs passed the physics criteria except for the 32-bit *fpzip* run.

The *fpzip* compression rate varied over the course of the simulation and at different regions within the simulation domain. Figure 12 shows the compression rate of all fields from a 1024 processor RTI simulation using Miranda and compressed with *fpzip* in 48-bit mode on each time step. The compression rate begins at ~ 2.2 in processors near the mixing layer and at ~ 4.5 in processors farther away. Over time the minimum compression rate decreases, reaching a floor in some processes at 1.7. Over the course of the simulation, an average compression rate of 3.3 was achieved.

Figure 13 shows the compression rate of one of the velocity components from a 1024 processor RTI simulation using Miranda and *APAX* 2X compression on each time-step. The average compression rate is slightly above 2X due to the fact that the *APAX* coder converts all float values into a 32-bit

signed integer representation before quantization. The minimum compression ratio never falls below 2.0 on any processor. A fixed rate compression scheme may be important for memory limited applications, in which a strict lower bound on compression rate is required.

6. DISCUSSION AND FUTURE WORK

Disk compression can have a positive impact on performance today. For example, the bandwidth per core to the parallel file system on Cielo, a Cray XE-6 at Los Alamos National Laboratory, is roughly 1 MB/s for pF3D when 1k to 64k cores write simultaneously. Lossy compression algorithms have a throughput of roughly 200 MB/s, so the time spent compressing data would be small compared to the time spent writing it to disk. Thus the total time spent writing to disk would be reduced by roughly the compression ratio. For checkpointing, file formats such as HDF5 already provide support for lossless compression, and *APAX* is available as a compression filter for HDF5. The flexibility of libraries make it straightforward to also add support for other lossy compression schemes.

Memory compression will not help performance on current multi-core systems because memory bandwidth per core (roughly 8.6 GB/s on an Intel Sandy Bridge) is significantly higher than software compression rates. An interesting possibility is to put compression in on-chip logic. The *APAX* compression algorithm can achieve throughputs of 2.5 to 8 GB/s while using between 0.1 and 0.4 mm² in 28 nm CMOS and increasing the power consumption of a chip by less than one percent. This rate is close enough to the memory bandwidth per core that memory compression might increase overall application performance. We expect that hardware compression schemes will be required to make “memory compression” a generally useful approach.

To take full advantage of memory compression, it would be necessary to not only compress and decompress the simulation state at the end of each cycle, but to also insert compression into all reads and writes to and from main memory, possibly storing localized views of uncompressed simulation state in scratch memory or cache. To realize this scenario, we envision having to partition the data arrays into bite-size compressed chunks that can be processed and evolved mostly independently. We see such a restructuring of inner loops in the simulation code, possibly assisted by semi-automated techniques, as future work. In particular, a separate study is needed to assess the impact of data blocking and more frequent compression on accuracy and speed.

Another avenue for future work is the development of models to predict the effects of lossy compression. In our current approach, the simulation is run repeatedly with increasing compression ratios and corresponding errors to determine the failure point below which compression is deemed “safe.” Analogous to how it is possible to analyze the impact mesh resolution has on accuracy, we require models to predict how quantization and variable-precision arithmetic influences the simulation state and error growth over time. The *APAX* Profiler [32] has been used successfully to predict acceptable compression rates for other media by analyzing the intrinsic noise level relative to the signal strength. However, such tools do not account for the potential cascading effects of compression-induced errors that grow over time, nor how these errors correlate with physics-based metrics.

7. CONCLUSIONS

The LULESH study indicates that a broad range of compression parameters are valid for this run. The Sedov blast wave problem is a stable problem, and we expect that for unstable flows higher precision settings may be required. We noted that *APAX* and *fpzip* were similar in their performance, indicating that compression schemes with either fixed rate or fixed precision modes are viable. We found that with *APAX*, the per-block scaling of values did lead to unphysical noise in the nodal positions in front of the shock wave at a requested compression rate of 4X. In future work, it would be interesting to thread the compression calls through the LULESH solver at the point when each field value is updated, as this would be closer to how memory compression would occur in practice.

The Miranda run with 32-bit *fpzip* failed when compressed at every time step but passed when run at the much lower frequency required for “disk compression.” This run had a compression rate of 4.5X, only slightly higher than the successful 40-bit *fpzip* and *APAX* 4X runs. We showed that *fpzip* generates high wavenumber features in the density spectrum, and that these interact with the numerics in Miranda to cause the problem to crash. With the *APAX* compressor, Miranda successfully ran with compression rates of 5.3X, and we showed that it did not exhibit the high frequency features. A key result is that the details of how compression algorithms quantize and round can be important factors in overall results when compression is integrated with numerical algorithms.

A pF3D run with 32-bit *fpzip* (average compression ratio 3.66X) had low errors while a 24-bit *fpzip* run (average compression ratio 6.57X) failed. The pF3D simulations using *APAX* passed the physics criteria at compression ratios up to 6X. These results show that the onset of failure occurs across a fairly narrow range of compression ratios and can depend on interactions between the compressor and the code’s numerics.

This paper has examined the impact of lossy compression on three physics simulation codes – LULESH, Miranda, and pF3D. All three codes can be run with a lossy compression ratio of 3X or greater when using a frequency suitable for “memory compression.” Lossy compression is also a viable approach to reducing the impact of limited disk bandwidth for all three codes. To our knowledge, this is the first study of the effects of applying lossy compression to the physics state of simulations as a strategy for mitigating the data movement bottleneck expected on future systems.

Acknowledgment

This document was released as LLNL-CONF-642236. This work was performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under Contract DE-AC52-07NA27344. Some of this material is based upon work supported by the Defense Advanced Research Projects Agency under its Agreement No. HR0011-07-9-0001.

8. REFERENCES

- [1] R. L. Berger, B. F. Lasinski, A. B. Langdon, T. B. Kaiser, B. B. Afeyan, B. I. Cohen, C. H. Still, and E. A. Williams. Influence of spatial and temporal laser beam smoothing on stimulated brillouin scattering in filamentary laser light. *Physical Review Letters*, 75(6):1078–1081, Aug. 1995.
- [2] M. Burtscher and P. Ratanaworabhan. High throughput compression of double-precision floating-point data. In *Data Compression Conference*, pages 293–302, Mar. 2007.
- [3] W. H. Cabot and A. W. Cook. Reynolds number effects on Rayleigh-Taylor instability with possible implications for type-Ia supernovae. *Nature Physics*, 2:562–568, 2006.
- [4] E. J. Candès and M. B. Wakin. An introduction to compressive sampling. *IEEE Signal Processing Magazine*, 25(2):21–30, Mar. 2008.
- [5] A. W. Cook, W. Cabot, and P. L. Miller. The mixing transition in Rayleigh–Taylor instability. *Journal of Fluid Mechanics*, 511:333–362, 2004.
- [6] A. W. Cook, W. H. Cabot, M. L. Welcome, P. L. Williams, B. J. Miller, B. R. de Supinski, and R. K. Yates. Tera-scalable algorithms for variable-density elliptic hydrodynamics with spectral accuracy. In *ACM/IEEE Conference on Supercomputing*, page 60, 2005.
- [7] D. L. Donoho. Compressed sensing. *IEEE Transactions on Information Theory*, 52(4):1289–1306, Apr. 2006.
- [8] N. Fout and K.-L. Ma. An adaptive prediction-based approach to lossless compression of floating-point volume data. *IEEE Transactions on Visualization and Computer Graphics*, 18(12):2295–2304, Dec. 2012.
- [9] N. Huebbe and J. Kunkel. Reducing the HPC-datastorage footprint with MAFISC multidimensional adaptive filtering improved scientific data compression. *Computer Science Research and Development Journal*, 28(2–3):231–239, May 2012.
- [10] N. Huebbe, A. Wegener, J. Kunkel, Y. Ling, and T. Ludwig. Evaluating lossy compression on climate data. In *International Supercomputing Conference*, pages 343–356, June 2013.
- [11] L. Ibarria, P. Lindstrom, J. Rossignac, and A. Szymczak. Out-of-core compression and decompression of large n -dimensional scalar fields. *Computer Graphics Forum*, 22(3):343–348, 2003.
- [12] I. Karlin, A. Bhatele, B. L. Chamberlain, J. Cohen, Z. Devito, M. Gokhale, R. Haque, R. Hornung, J. Keasler, D. Laney, E. Luke, S. Lloyd, J. McGraw, R. Neely, D. Richards, M. Schulz, C. H. Still, F. Wang, and D. Wong. LULESH programming model and performance ports overview. Technical Report LLNL-TR-608824, Lawrence Livermore National Laboratory, Dec. 2012.
- [13] I. Karlin, A. Bhatele, J. Keasler, B. L. Chamberlain, J. Cohen, Z. DeVito, R. Haque, D. Laney, E. Luke, F. Wang, D. Richards, M. Schulz, and C. Still. Exploring traditional and emerging parallel programming models using a proxy application. In *IEEE International Parallel & Distributed Processing Symposium*, pages 919–932, May 2013.

- [14] J. Keasler and R. Hornung. Hydrodynamics challenge problem. Technical Report LLNL-TR-490254, Lawrence Livermore National Laboratory, 2010.
- [15] P. Kogge, K. Bergman, S. Borkar, D. Campbell, et al. Exascale computing study: Technology challenges in achieving exascale systems, Sept. 2008.
- [16] S. Lakshminarasimhan, N. Shah, S. Ethier, S. Klasky, R. Latham, R. Ross, and N. F. Samatova. Compressing the incompressible with ISABELA: In-situ reduction of spatio-temporal data. In *Euro-Par Parallel Processing, Lecture Notes in Computer Science*, pages 366–379. Springer, 2011.
- [17] S. Langer, B. Still, T. Bremer, D. Hinkel, B. Langdon, J. A. Levine, and E. A. Williams. Cielo full-system simulations of multi-beam laser-plasma interaction in NIF experiments. In *Proceedings of the 53rd Cray User Group Meeting*, 2011.
- [18] H. Lehmann and B. Jung. In-situ data compression for flow simulation in porous media. In *Parallel & Distributed Processing Techniques & Applications*, 2012.
- [19] P. Lindstrom. fpzip version 1.0.1, 2008. <https://computation.llnl.gov/casc/fpzip/>.
- [20] P. Lindstrom and M. Isenburg. Fast and efficient compression of floating-point data. *IEEE Transactions on Visualization and Computer Graphics*, 12(5):1245–1250, 2006.
- [21] S. K. Moore. Multicore is bad news for supercomputers. *IEEE Spectrum*, 45(11):15, 2008.
- [22] E. I. Moses. Overview of the National Ignition Facility. *Fusion Science and Technology*, 54(2):361–366, 2008.
- [23] E. I. Moses, R. N. Boyd, B. A. Remington, C. J. Keane, and R. Al-Ayat. The National Ignition Facility: Ushering in a new age for high energy density science. *Physics of Plasmas*, 16(041006):1–13, 2009.
- [24] S. Muraki. Approximation and rendering of volume data using wavelet transforms. In *IEEE Visualization*, pages 21–28, 1992.
- [25] S. Muraki. Multiscale volume representation by a dog wavelet. *IEEE Transactions on Visualization and Computer Graphics*, 1(2):109–116, 1995.
- [26] R. Murphy. On the effects of memory latency and bandwidth on supercomputer application performance. *IEEE International Symposium on Workload Characterization*, pages 34–43, 2007.
- [27] E. Schendel, Y. Jin, N. Shah, J. Chen, C. S. Chang, S.-H. Ku, S. Ethier, S. Klasky, R. Latham, R. Ross, and N. Samatova. ISOBAR preconditioner for effective and high-throughput lossless data compression. In *IEEE International Conference on Data Engineering*, pages 138–149, 2012.
- [28] C. H. Still, R. L. Berger, A. B. Langdon, D. E. Hinkel, L. J. Suter, and E. A. Williams. Filamentation and forward Brillouin scatter of entire smoothed and aberrated laser beams. *Physics of Plasmas*, 7(5):2023–2032, 2000.
- [29] E. Tasker, R. Brunino, N. Mitchell, D. Michielsen, S. Hopton, F. Pearce, G. Bryan, and T. Theuns. A test suite for quantitative comparison of hydrodynamics codes in astrophysics. *Monthly Notices of the Royal Astronomical Society*, 390(3):1267–1281, 2008.
- [30] A. Wegener. Adaptive compression and decompression of bandlimited signals, Mar. 2006. http://www.patentlens.net/patentlens/patent/US_7009533/.
- [31] A. Wegener. Block floating point compression of signal data, Oct. 2012. http://www.patentlens.net/patentlens/patent/US_8301803/.
- [32] A. Wegener. Universal numerical encoder and profiler reduces computing memory wall with software, FPGA, and SoC implementations. In *IEEE Data Compression Conference*, page 528, Mar. 2013.
- [33] A. Wegener, N. Chandra, Y. Ling, R. Senzig, and R. Herfkens. Effects of fixed-rate CT projection data compression on perceived and measured CT image quality. In *SPIE Medical Imaging*, volume 7627, Feb. 2010.
- [34] M. L. Wilkins. *Methods in Computational Physics*. Academic Press, 1964.
- [35] J. Woodring, S. Mniszewski, C. Brislawn, D. DeMarle, and J. Ahrens. Revisiting wavelet compression for large-scale climate data using JPEG 2000 and ensuring data precision. In *IEEE Large Data Analysis and Visualization*, pages 31–38, 2011.