

High Performance Multivariate Visual Data Exploration for Extremely Large Data

Oliver Rübél^{1,5,6}, Prabhat¹, Kesheng Wu¹, Hank Childs², Jeremy Meredith³, Cameron G.R. Geddes⁴,
Estelle Cormier-Michel⁴, Sean Ahern³, Gunther H. Weber¹, Peter Messmer⁷, Hans Hagen⁶,
Bernd Hamann^{1,5,6} and E. Wes Bethel^{1,5}

Abstract—One of the central challenges in modern science is the need to quickly derive knowledge and understanding from large, complex collections of data. We present a new approach that deals with this challenge by combining and extending techniques from high performance visual data analysis and scientific data management. This approach is demonstrated within the context of gaining insight from complex, time-varying datasets produced by a laser wakefield accelerator simulation. Our approach leverages histogram-based parallel coordinates for both visual information display as well as a vehicle for guiding a data mining operation. Data extraction and subsetting are implemented with state-of-the-art index/query technology. This approach, while applied here to accelerator science, is generally applicable to a broad set of science applications, and is implemented in a production-quality visual data analysis infrastructure. We conduct a detailed performance analysis and demonstrate good scalability on a distributed memory Cray XT4 system.

I. INTRODUCTION

This work focuses on combining and extending two different but complementary technologies aimed at enabling rapid, interactive visual data exploration and analysis of contemporary scientific data. To support highly effective visual data exploration, knowledge discovery and hypothesis testing, we have adapted and extended the concept of parallel coordinates, in particular binned or histogram-based parallel coordinates, for use with high performance query-driven visualization of very large data. In the context of visual data exploration and hypothesis testing, the parallel coordinates display and interaction mechanism serves multiple purposes. First, it acts as a vehicle for visual information display. Second, it serves as the basis for the interactive construction of compound Boolean

data range queries. These queries form the basis for subsequent “drill down” or data mining actions. To accelerate data mining, we leverage state-of-the-art index/query technology to quickly mine for data of interest as well as to quickly generate multiresolution histograms used as the basis for the visual display of information. This combination provides the ability for rapid, multiresolution visual data exploration.

We apply this new technique to large, complex scientific data created by a numerical simulation of a laser wakefield particle accelerator. In laser wakefield accelerators, particles are accelerated to relativistic speeds upon being “trapped” by the electric fields of plasma density waves generated by the radiation pressure of an intense laser pulse fired into the plasma. These devices are of interest because they are able to achieve very high particle energies within a relatively short amount of distance when compared to traditional electromagnetic accelerators. The VORPAL [1] simulation code is used to model experiments, such as those performed at the LOASIS facility at LBNL [2], and is useful in helping to gain deeper understanding of phenomena observed in experiments, as well as to help formulate and optimize the methodology for future experiments.

Laser wakefield simulations model the behavior of individual particles as well as the behavior of the plasma electric and magnetic fields. Output from these simulations can become quite large: today’s datasets, such as the ones we study here, can grow to be on the order of 200GB per timestep, with the simulation producing ≈ 100 timesteps. The scientific challenge we help address in this study is first to quickly find particles that have undergone wakefield acceleration, then trace them through time to understand acceleration dynamics, and perform both visual and quantitative analysis on the set of accelerated particles.

One scientific impact of our work is that we have vastly reduced the duty cycle in visual data exploration and mining. In the past, accelerator scientists would perform the “trace backwards” step using scripts that performed a search at each timestep for a set of particles. Runtimes for this operation were on the order of hours. Using our implementation, those runtimes are reduced from hours to seconds.

The specific new contributions of this work are as follows:

- We present a novel approach for quickly creating histogram-based parallel coordinates displays. These displays serve to convey information as well as the interface for the interactive construction of compound, multivariate

1. Computational Research Division, Lawrence Berkeley National Laboratory, One Cyclotron Road, Berkeley, CA 94720, USA.

2. Lawrence Livermore National Laboratory, 7000 East Avenue, Livermore, CA 94550.

3. Oak Ridge National Laboratory, P.O. Box 2008, Oak Ridge, TN 37831, USA.

4. LOASIS program of Lawrence Berkeley National Laboratory, One Cyclotron Road, Berkeley, CA 94720, USA.

5. Institute for Data Analysis and Visualization (IDAV) and Department of Computer Science, University of California, Davis, One Shields Avenue, Davis, CA 95616, USA.

6. International Research Training Group “Visualization of Large and Unstructured Data Sets – Applications in Geospatial Planning, Modeling, and Engineering,” Technische Universität Kaiserslautern, Erwin-Schrödinger-Straße, D-67653 Kaiserslautern, Germany.

7. Tech-X Corporation, 5621 Arapahoe Ave. Suite A, Boulder, CO 80303.

Boolean range queries. The new approach leverages state-of-the-art index/query technology to achieve very favorable performance rates.

- We apply this new approach to solve a challenging scientific data understanding problem in accelerator modeling. This new approach performs particle tracking in a few seconds as compared to hours when using a naive script.
- We examine and report the performance of our approach using a modern HPC platform on a very large ($\approx 1.5\text{TB}$) and complex scientific dataset. We demonstrate that our approach has excellent scalability characteristics on a distributed memory Cray XT4 system.

The rest of this paper is organized as follows. First, we review relevant background work in Section II, which covers a diverse set of topics in visual information display, index/query technology, and high performance visual data analysis software architectures. Next in Section III, we present the architecture and implementation of our approach. This approach is then applied to solve a challenging scientific data understanding problem in the field of laser wakefield accelerator modeling in Section IV. We evaluate the performance of our system and present performance results in Section V. Finally, we conclude in Section VI with a summary and suggestions for future work.

II. RELATED WORK

A. Information Display – Parallel Coordinates

Parallel coordinates, proposed by Inselberg [3] and Wegman [4], are a common information visualization technique for high-dimensional data sets. In parallel coordinates, each data variable of a multivariate dataset is represented by one axis. The parallel coordinates plot is constructed by drawing a polyline connecting the points where a data record's variable values intersect each axis. This type of plot is expensive in the sense that many pixels are required to represent a single data record. In this form of display, there is substantial data occlusion by the many polylines required to display all records from a large dataset. An overview of modern parallel coordinates is provided in [5], [6].

Fua et al. [7] proposed using hierarchical parallel coordinates, based on hierarchical clustering, to create a multiresolution view of the data that enables data exploration at varying levels of detail. Johannson et al. [8] used clustering to determine the inherent structure of data, and displayed that structure with high-precision textures using different texture transfer functions. Novotný then used a binning algorithm based on a *k-means* clustering approach for creating an aggregate parallel coordinates visualization [9]. All these approaches are well-suited for presenting static data, but are not well suited for time-varying data since defining temporally consistent clusters is non-trivial and computationally expensive.

Our histogram-based parallel coordinates approach extends the work of Novotný and Hauser [10], who proposed using binned parallel coordinates as an output-oriented (rendering) approach. The main limitation of this approach is that in

order to achieve interactive rendering speed, their method precomputes all possible 2D histograms with a fixed resolution of 256×256 and uniform, equal-sized histogram bins. To implement different level-of-detail views, bins in the pre-computed histograms are merged, reducing the number of bins by half in each drill down step. While this strategy is efficient, it has several limitations. First, it does not allow for smoothly drilling into finer-resolution views of the data. Second, it supports presentation of binned views of the entire dataset, but does not support user-defined data subsetting. Finally, the fixed 256×256 histogram resolution exhibits significant aliasing when zooming in on narrow variable ranges in the parallel coordinates plot. Their approach uses binned parallel coordinates for “context views” and traditional, polyline-based parallel coordinates for “focus views.” These focus views may still contain a substantial number of data records and will suffer from extensive occlusion as a result.

B. High Performance Index/Query for Data Mining

The data access patterns for data mining and analysis applications tend to be markedly different from those for transaction-based applications. Transaction-based applications tend to read and then update records in a database. In contrast, data mining and analysis applications tend to be read-only in their access pattern. The database indexing technology that is best suited for this type of data access is known as the bitmap index [11], [12]. The core idea of a bitmap index is to use a sequence of bits to mark the positions of records satisfying certain conditions. Searching bitmap indices for data records that match a set of range conditions is performed efficiently using Boolean operations on bit vectors.

In uncompressed form, such bitmap indices may require too much space for variables with many distinct values, such as particle position or momentum. Several techniques have been proposed to improve the efficiency of the bitmap indexes on such variables [13], [14]. We make use of a bitmap index software called FastBit [15]. It implements the fastest known bitmap compression technique [16], [17], and has been demonstrated to be effective in a number of data analysis applications [18], [19]. In particular, it has a number of efficient functions for computing conditional histograms [20], which are crucial for this work. Furthermore, FastBit indices are relatively small compared to popular indices such as B-trees [16, Fig. 7] and can be constructed much faster than others [21, Fig. 12]. Bitmap indices are well-known for their effectiveness on data with relatively small number of distinct values, such as gender. FastBit indices have been demonstrated to be very efficient also for data with a large number of distinct values through its unique compression [14] and binning [22].

When a variable has a large number of distinct values the corresponding FastBit index is typically on a binned version of the data, where a bit is 1 if the value of a record falls in a particular bin. In this case, the number of 1s in a bitmap corresponding to a bin is the number of records in the bin. This provides an efficient method for computing a histogram [20]. FastBit offers a number of different options

for creating bitmap bins. When composing range queries, users typically specify conditions with relatively low-precision values, such as pressure less than $1 * 10^{-5}$ or momentum greater than $2.5 * 10^8$. The constant $1 * 10^{-5}$ is said to have 1-digit precision, and the constant $2.5 * 10^8$ to have 2-digit precision. FastBit can build indices with bin boundaries with any user-specified precision so that all queries involving low-precision boundaries are answered accurately with index only. These features make FastBit uniquely suitable for this work.

C. High Performance Query-Driven Visualization

The combination of high performance index/query with visual data exploration tools was described by Stockinger et al. [23] using the term “Query-driven visualization.” That work focuses on comparing the performance of such a combination with state-of-the-art, tree-based searching structures that form the basis for a widely-used isocontouring implementation. Their work shows that this approach outperforms tree-based search structures for scalar variables, and also points out that all tree-based index/search structures are not practical for large, multivariate datasets since they suffer from the “Curse of Dimensionality” [24]. The basic idea there is that storage complexity grows exponentially as one adds more and more search dimensions (e.g., more variables to be indexed/searched).

These concepts were later extended to the analysis of massive collections of network traffic data in two related works. First, the notion of performing network traffic analysis using statistics (e.g., histograms) rather than raw data led to a methodology that enabled exploration and data mining at unprecedented speed [18]. That study showed using these concepts to rapidly detect a distributed scan attack on a dataset of unprecedented size – 2.5 billion records. There, users are presented with an interface consisting of histograms of individual variables, and then they formulate a complex query via a process that is essentially a histogram “cross product.” The process of data mining was subsequently accelerated through a family of algorithms for computing conditional histograms on SMP parallel machines [25].

D. High Performance Visual Data Analysis

Childs et al. [26] demonstrated a data processing and visualization architecture that is capable of scaling to extreme dataset sizes. This software system, VisIt [27], has been shown to scale to tens of billions of data points per timestep and runs in parallel on nearly all modern HPC platforms. In reference to our work, VisIt employs a contract-based communication system that allows for pipelining and I/O optimization, reducing unnecessary processing and disk access. It is this extensible contract system that we utilize to generate histograms for visual data exploration. A specific optimization that we employ is VisIt’s concept of a set of Boolean range queries. This set is communicated between data processing modules in an out-of-band fashion, allowing downstream filters to limit the scope of work of upstream filters.

VisIt’s remote visualization capability crosses several independent axes. In addition to parallelization of data fetching,

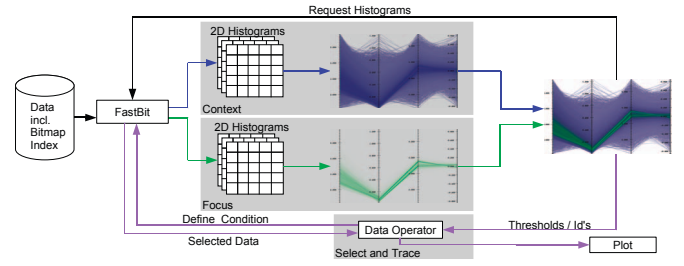


Fig. 1. An overview of the major components and data flow paths in our implementation. Large-scale scientific data and indexing metadata are input via a parallel I/O layer, which allows us to achieve high levels of performance and parallel efficiency. We use an API at the I/O layer to perform parallel computation of multidimensional histograms as well as data subsetting. Results of those computations are used downstream in the visualization application for presenting information to the user and in support of interactive data mining actions.

VisIt also performs data extraction and calculation entirely in parallel. Finally, rendering may be done in serial or parallel, depending on the data load. If the resultant geometry is small enough, it is collected at the HPC side and shipped across the network to the user’s desktop for GPU-based rendering. However, if the geometry is too large for interactive display on a single GPU, VisIt employs sort-last rendering and compositing on the HPC system, with the resultant pixels shipped across the network to the user’s display. The model described here has many parallels to remote visualization architectures such as those in ParaView [28] and EnSight [29].

III. SYSTEM DESIGN

Figure 1 shows a high-level view of the components and data flow in our implementation. Raw scientific data, which is produced by simulation or experiment, is augmented by the computation of indexing data. In our case, this step is performed outside the visual data analysis application as a one-time preprocessing, and our implementation uses FastBit [15] for creating index structures. The data sizes are described in more detail in Section IV and V.

After the one-time preprocessing step, our implementation uses FastBit at the data-loading portion of the pipeline to quickly compute histograms in parallel and to perform high-performance data subsetting/selection based upon multivariate thresholds or particle identifier (Section III-B). These histograms serve as the basis for the visual presentation of full-resolution and subset views of data vis-a-vis parallel coordinates plots (Section III-A). In our implementation, the computational complexity of rendering parallel coordinates plots – both context and focus views – is a function of histogram resolution, not the size of the underlying data. Therefore, our approach is particularly well-suited for application to extremely large data, which we present in Section IV.

A. Histogram based Parallel Coordinates

Parallel coordinates provide a very effective interface for defining multi-dimensional queries based on thresholding. Using sliders attached to each axis of the parallel coordinates

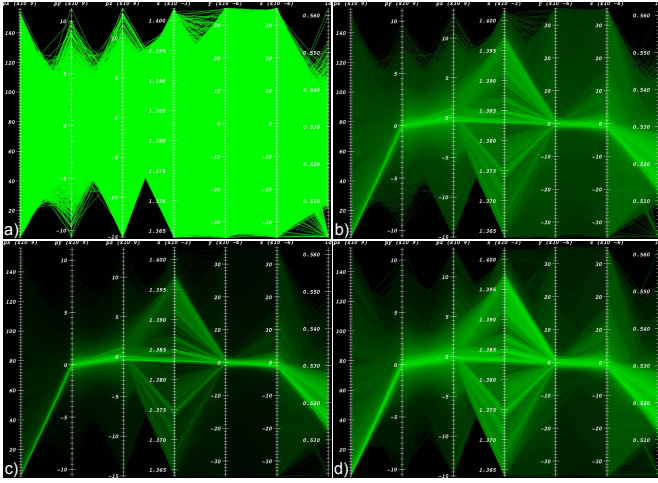


Fig. 2. Comparison of different parallel coordinate renderings of a subset of a 3D laser wakefield particle acceleration dataset consisting of 256,463 data records and 7 data dimensions. a) Traditional line based parallel coordinates. b) High-resolution, histogram-based parallel coordinates with 700 bins per data dimension. c) Same as in b, but using a lower gamma value g defining the basic brightness of bins. d) Same rendering as in b but using only 80 bins per data dimension. When comparing a and b, we see that the histogram-based rendering reveals many more details when dealing with a large number of data records. As illustrated in c, by lowering gamma we can then reduce the brightness of the plot and even remove sparse bins, thereby producing a plot that focuses on the main, dense features of the data. By varying the number of bins we can then create renderings at different levels of detail.

plot, a user defines range thresholds in each displayed dimension. By rendering the user-selected data subset (the focus view) in front of a parallel coordinates plot created from the entire data set (or in many cases a subset of it) (the context view), the user receives immediate feedback about general properties of the selection. Data outliers stand out visually as single or small groups of lines diverging from the main data trends. Data trends appear as dense groups of lines (bright colored bins in our case). A quick visual comparison of the focus and context views helps to convey understanding about similarities and differences between the two.

In practice, parallel coordinates have disadvantages when applied to very large datasets. First, each data record is represented with a single polyline that connects each of the parallel coordinates axes. As data size increases, the plot becomes more cluttered and difficult to interpret. Also, data records drawn later will occlude information provided by data records drawn earlier. Worst of all is that this approach has computational and rendering complexity that is proportional to the size of the dataset. As data sizes grow ever larger, these problems become intractable.

To address these problems, we employ an efficient rendering technique based on two-dimensional (2D) histograms. Rather than viewing the parallel coordinates plot as a collection of polylines, one per data record, we approach rendering by considering instead the relationships of all data records between pairs of parallel coordinate axes. That relationship can be discretized as a 2D histogram and then later rendered. This idea was introduced in earlier work [10]. As illustrated in Figure 3, we create a parallel coordinates representation based on 2D

histograms by drawing one quadrilateral per non-empty bin, where each quadrilateral connects two data ranges between neighboring axes. As illustrated in Figure 2(a,b), histogram-based rendering overcomes the limitations of polyline-based rendering and reveals much more data detail when dealing with a large number of data records.

In the following, we first describe how we quickly compute 2D histograms, and then explain in detail how we use them to efficiently render parallel coordinates. Having introduced the principle of histogram-based parallel coordinates, we then compare use of uniformly (equal width) and adaptively (equal weight) binned 2D histograms in the context of rendering parallel coordinate views.

1) *Computing Histograms:* In this work, we implement the computation of 2D histograms at the data I/O stage of VisIt: 2D histograms are computed directly in the file reader, which leverages FastBit for index/query operations as well as histogram computation. This approach has several major benefits within the context of very large, high performance visual data analysis. First, instead of having to read the entire data set and transfer it to the plot to create a rendered image, we limit internal data transfer and processing to a set of 2D histograms, which are very small when compared to the size of the source data. Second, data I/O is limited to those portions needed for computing the current 2D histogram, i.e., information on the relevant particles in two data dimensions. After computing a 2D histogram, we can then discard the data the loader may have read to compute the histogram, thus decreasing the memory footprint. An added advantage is the fact that all the computationally and I/O intensive work is done at the beginning of the parallel execution pipeline rather than at the end in the plot. Third, having the histogram computation be part of the file reader allows us to perform such computation at the same stage of processing as parallel I/O, which is one of the most expensive operations in visual data analysis. This approach provides the ability to achieve excellent parallel performance, as described in Section V.

2) *Using Histograms to Render Parallel Coordinates Plots:* When using 2D histograms for rendering parallel coordinates plots, we have access to additional information that is not present when using traditional polyline-based methods. We know, for example, the number of records contributing to specific variable ranges between parallel axes. This extra information allows us to optimize various aspects of data visualization to convey more information to the user. We use brightness, for example, to reflect the number of records per bin, which leads to improved visual presentation. Assuming that denser regions are more important than sparse regions, we render bins in back-to-front order with respect to the number of records per bin $h(i, j)$. Figure 2(a,b) shows a direct comparison of the same data once rendered using traditional line based parallel coordinates and once using our histogram based rendering approach.

In order to further improve the rendering, we allow the user to define a gamma value g defining the overall brightness of the plot. As illustrated in Figure 2c, a lower g value will reduce

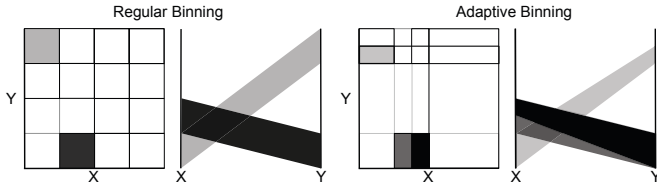


Fig. 3. Illustration showing rendering of parallel coordinates based on uniformly (left) and adaptively binned histograms (right). By using higher resolution in areas of extremely high density, an adaptive binning is able to represent the general data trends much more accurately.

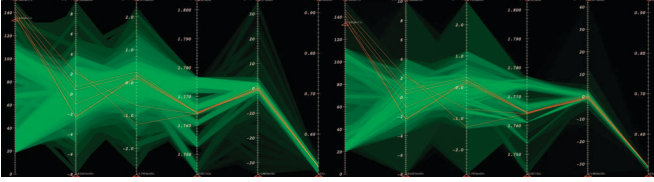


Fig. 4. Histogram-based parallel coordinates using 32×32 uniformly binned histograms (left) and adaptive histograms (right). Compared to uniform binning, the adaptive binning preserves more details in dense areas while discarding some details in sparse areas of the data. An adaptively binned plot may ease comparison of selections (red) with dense structures present in the data in low level of detail views.

the brightness of the plot, or even remove sparse bins from the rendering, thereby producing a much less cluttered visualization that focuses attention on the main, dense data features. Since our method is not constrained by a fixed histogram bin resolution, we can easily recompute histograms at higher resolution, or, using adaptive binning, produce visualizations at varying levels of detail. This feature is important for providing smooth drill-down into finer levels of detail in very large datasets and represents one of the major improvements of our approach over previous work. Figure 2d, shows as an example the same rendering as in Figure 2b, but using just 80 bins per data dimension. Another example is also provided in Figure 4.

Previous histogram-based parallel coordinates work used histogram-based rendering for the context view and traditional line-based rendering for the focus view. One limitation of this approach is that the focus view may consist of a very large number of data records. We overcome this limitation by using a histogram-based approach for both the context and focus views. This approach is feasible given the rapid rate at which we can recompute new conditional 2D histograms as explained in Section V. The focus view is rendered on top of the context view using a different color to make the focus more easily distinguishable. This approach has the further advantage that we can render it at different levels of detail simply by specifying the number of bins per variable. Here, we can use low-resolution histograms for the context view and a high-resolution histograms for the focus view, thereby supporting a high quality and smooth visual drill-down into the focus view.

Similar to rendering the focus on top of the context, we can also create a rendering of multiple timesteps in one parallel coordinates plot. To do so, we assign a unique color to each timestep and render the individual plots, each representing one timestep, on top of one other. As we will show later in Section IV, such a rendering can be helpful to identify the

general temporal changes in the data. In practice, temporal parallel coordinates are most useful when analyzing some characteristic subset of the data.

3) *Adaptive and Uniform Histogram Binning*: Regular histograms, with uniform, equal-sized bins, are well-suited for high-resolution renderings but have substantial disadvantages when creating low-level-of-detail views where the number of bins per variable is much smaller than the number of pixels per parallel axis. In such a case each bin may cover regions of varying data density, e.g., large areas of a bin may contain only few data records while other areas contain many. Our approach offers another improvement over previous work, namely the possibility to compute and render with adaptive, rather than uniform histogram bins. With adaptive binning, each bin of the histogram contains approximately the same number of data records, which may offer advantages in certain situations [30].

As an example, Figure 4 shows a comparison of data visualized using 32×32 uniform versus adaptive histogram bins. In comparison to a uniform binning, adaptive binning discards some features in sparse areas of the data to preserve more information in dense areas. Adaptively binned histograms may ease comparison of selections with general data trends. As illustrated in Figure 3, when using adaptively binned histograms, a more generalized rendering is required to allow rectangles to connect different-size ranges on neighboring axes. Also, since the area $a(i, j)$ covered by each bin is no longer constant, we need to compute the brightness of each bin, and assign rendering order based on the actual data density per bin $p(i, j) = \frac{h(i, j)}{a(i, j)}$, rather than directly based on the number of records per bin $h(i, j)$. Uniform binning is well-suited for high-resolution renderings of the data while adaptive binning may be advantageous for low-resolution renderings.

In other applications, such as statistical analysis, one is often interested in both the main data trends as well as outlier behavior. In order to achieve an optimal low-level-of-detail rendering for such applications, one could further restrict the minimal density p of bins during computation of the adaptive binning to ensure that details in sparse areas of the data are represented accurately. As proposed by Novotný and Hauser [10] one may also employ a separate outlier detection scheme for this purpose in which data records located in bins of extremely low density p are rendered as individual lines resulting in a hybrid approach of line-based and histogram-based parallel coordinates.

B. Data Selection

In our implementation, there is synergy between presentation of information to the user and specification of queries. The parallel coordinates plot serves to present context and focus data views to the user, and also serves as the mechanism for specifying a multivariate Boolean range query. In our example case study in Section IV, a multivariate range query might take the form of $px > 10^9 \ \&\& \ py < 10^8 \ \&\& \ y > 0$, which selects high momentum particles in the upper half of the beam. Similar conditions can be formulated to define arbitrary particle-subsets with interesting momentum and spatial characteristics.

Once the user specifies such a multivariate range condition, those conditions are passed back upstream in the system to the FastBit-enhanced HDF5 reader for processing, either to compute new histograms or to extract data subsets that match the query for downstream processing. In the case of data subsetting, FastBit will locate those data records that satisfy the query and then pass them along for downstream processing. For conditional histograms, FastBit will compute new histograms using the query conditions as well as a histogram specification: the number of bins and the bin boundaries.

Once an interesting subset has been identified, yet another form of query can be issued in order to identify the same data subset (here particles) at different points in time. This type of query is of the form *ID IN* (id_1, id_2, \dots, id_n), where there are n particles in the subset. Again, such queries can be processed efficiently by FastBit and only the relevant set of particles is extracted and then passed along to visual data analysis machinery. This processing step offers a huge performance advantage: a technique without access to the index information must search the entire dataset for particle identifier matches. By issuing an identifier query across the entire time sequence, we construct particle tracks, which reveal valuable information about how particles are accelerated over time.

IV. USE CASE

We now present a specific example where we apply our system to perform visual data analysis of 2D and 3D data produced by a laser wakefield particle accelerator simulation. These simulations model the effects of a laser pulse propagating through a hydrogen plasma. Similar to the wake of a boat, the radiation pressure of the laser pulse displaces the electrons in the plasma, and with the space-charge restoring force of the ions, this displacement drives a wave (wake) in the plasma. Electrons can be trapped and accelerated by the longitudinal field of the wake, forming electron bunches of high energy. Due to the large amount of particles required in order to achieve accurate simulation results, it is not possible to simulate the entire plasma at once. As a result, the simulation is restricted to a window that covers only a subset of the plasma in x direction in the vicinity of the beam. The simulation code moves the window along the local x axis over the course of the run.

The simulation data in these examples contains information about the position of particles in physical space $-x, y, z$ –, particle momentum $-px, py, pz$ – and particle *ID* at multiple timesteps. As a derived quantity, we also include the relative particle position in x direction within the simulation window $x_{rel}(t) = x(t) - \max(x(t))$. There are more than 400,000 and 90 million particles per timestep in the 2D and 3D datasets, respectively. The 2D data consists of 38 timesteps and has an overall size of about 1.3GB, including the index structures. The 3D dataset consists of 30 timesteps and has an overall size of about 210GB, including the index. Each 3D timestep has a size of about 7GB, including the index, and about 5GB without the index. We perform a detailed analysis of the 2D dataset and then extend our analysis to the larger 3D dataset.

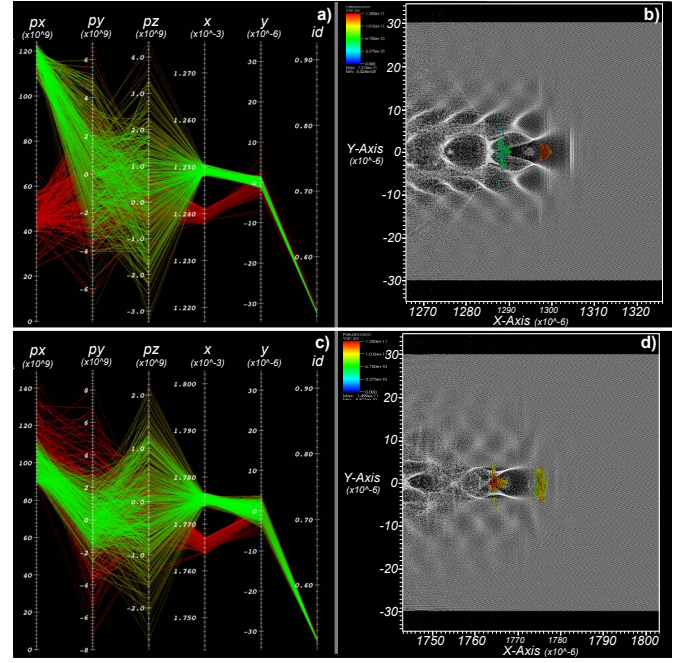


Fig. 5. a) Parallel coordinates and b) pseudocolor plot of the beam at $t = 27$. Corresponding plots c,d) at $t = 37$. The context plot, shown in red, shows both beams selected by the user after applying a threshold of $px > 8.872 \times 10^{10}$ at $t = 37$. The focus plot, shown in green, indicates the first beam that is following the laser pulse. In the pseudocolor plots b) and d), we show all particles in gray and the selected beams using spheres colored according to the particle's x -momentum, px . The focus beam is the rightmost bunch in these images. At timestep $t = 27$, the particles of the first beam (green in figure a) show much higher acceleration and a much lower energy spread (indicated via px) than the particles of the second beam. At later times, the lower momentum of the first beam indicates it has outrun the wave and moved into decelerating phase, e.g. at timestep $t = 37$.

In order to gain a deeper understanding of the acceleration process, we need to address complex questions such as: i) which particles become accelerated; ii) how are particles accelerated, and iii) how was the beam of highly accelerated particles formed and how did it evolve [31]. To identify those particles that were accelerated, we first perform selection of particles at a late timestep ($t = 37$) of the simulation by using a threshold for the value for x -momentum, px (Section IV-A). By tracing the selected particles over time we will then analyze the behavior of the beam during late timesteps (Section IV-B). Having defined the beam, we can analyze formation and evolution of the beam by tracing particles further back to the time where they entered the simulation and became injected into the beam (Section IV-C to IV-E). In this context, we use selection of particles, ID-based tracing of particles over time, and refinement of particle selections based on information from different timesteps as main analysis techniques.

As we will illustrate in this use case, the ability to perform selection interactively and immediately validate selections directly in parallel coordinates as well as other types of plots, such as pseudocolor or scatter-plots, enables much more accurate selection than previously possible. Tracing of particles over time then enables researchers to better understand evolution of the beam. With our visualization system the

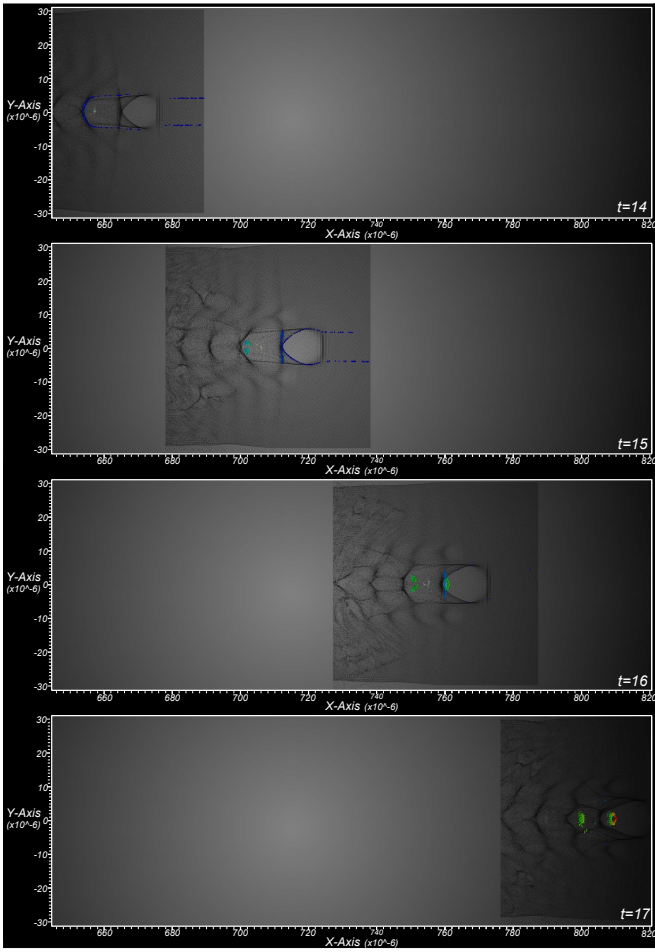


Fig. 6. Particles of the beam at timestep 14 (top left) to 17 (bottom right). The color of selected particles indicates x -momentum, px , using a rainbow color map (red=high, blue=low). Non-selected particles are shown in gray. We can readily identify two main sets of particles entering the simulation at timestep $t = 14$ and additional sets of particles entering at $t = 15$.

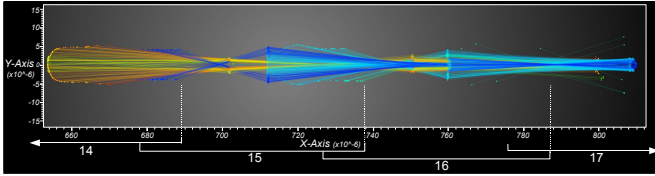


Fig. 7. The traces of the beam particles during timesteps $t = 14$ to $t = 17$ as shown in Figure 6. The positions of the selected particles at the different times are also shown. The begin and end of the timesteps in x direction are indicated below the bounding box. Here, we use color to indicate particle ID .

user can rapidly identify subsets of particles of interest and analyze their temporal behavior. For “small” datasets one will usually perform all analysis on a regular workstation. For larger datasets VisIt can also be run in a distributed fashion so that all heavy computation is performed on a remote machine where the data is stored and only the viewer and the GUI run on a local workstation.

A. Beam Selection

In order to identify the beam, i.e., find those particles that became accelerated, we first concentrate on the last timestep of

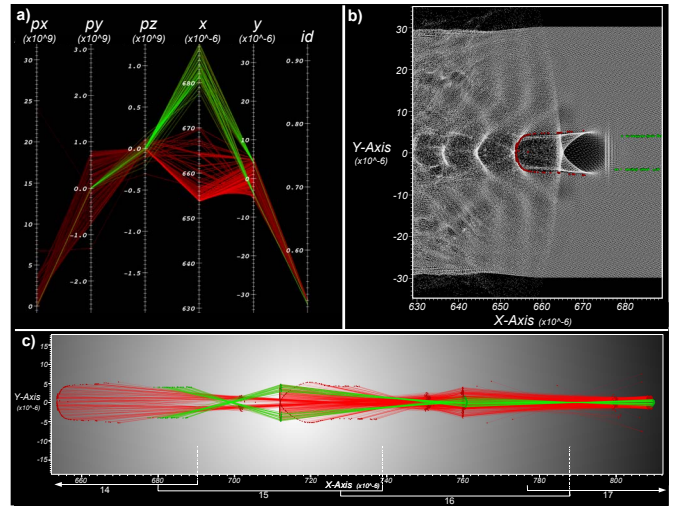


Fig. 8. a) By applying an additional threshold in x at timestep $t = 14$, we separate the two different set of particles entering the simulation. b) The refinement result, shown in physical space, includes all non-selected particles (gray) to provide context. c) Particle traces of the complete beam and the refined selection. In all plots we show the complete beam in red and the refined selection in green. After entering the simulation, the selected particles (green) define first the outer part of the first beam at timestep $t = 15$. Later on at timesteps $t = 16$ and $t = 17$, these particles become highly focused and define the center of the first beam.

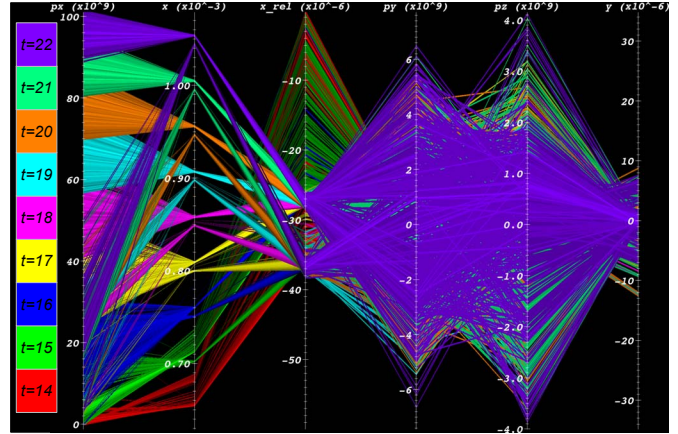


Fig. 9. The beam at timesteps $t = 14$ to $t = 22$ in a temporal parallel coordinates plots. Here, color indicates each of the discrete timesteps. We can readily identify the two different beams in x and x_{rel} ($x_{rel}(t) = x(t) - \max(x(t))$). While the second beam shows equal to higher values in px during early timesteps ($t = 14$ to $t = 17$), the first beam shows much higher acceleration at later times compared to the second beam.

the simulation (at $t = 37$). Using the parallel coordinates display, we select the particles of interest by applying a threshold of $px > 8.872 \times 10^{10}$. As is visible in the parallel coordinates plot, those particles constitute two separate clusters (beams) in x direction (Figure 5 c). Using a pseudocolor plot of the data, we can then see the physical structure of the beam.

B. Beam Assessment

By tracing particles back in time, we observe that the first bunch following the laser pulse (rightmost in these plots) has lower momentum spread at its peak energy (at $t = 27$) than the second bunch (see Figure 5a and 5b). In practice, the first beam

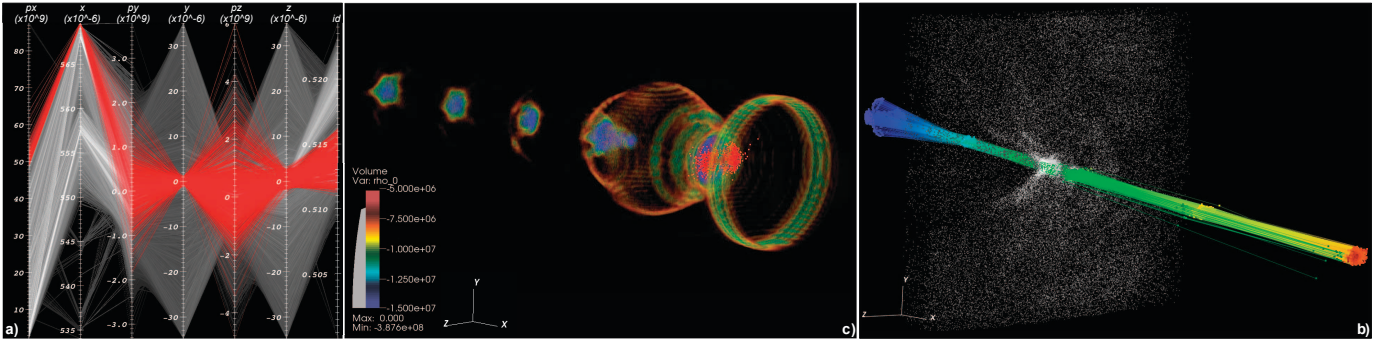


Fig. 10. a) Parallel coordinates of timestep $t = 12$ of the 3D dataset. Context view (gray) shows particles selected with $px > 2 \times 10^9$. The focus view (red) shows particles satisfying the condition $(px > 4.856 \times 10^{10}) \&\& (x > 5.649 \times 10^{-4})$, which form a compact beam in the first wake period following the laser pulse. b) Volume rendering of the plasma density and the selected focus particles (red). c) Traces of the beam. We selected particles at timestep $t = 12$, then traced the particles back in time to timestep $t = 9$ when most of the selected particles entered the simulation window. We also traced the particles forward in time to timestep $t = 14$. Color indicates px . In addition to the traces and the position of the particles, we also show the context particles at timestep $t = 12$ in gray to illustrate where the original selection was performed. We can see that the selected particles are constantly accelerated over time (increase in px).

following the laser pulse is therefore typically the one of most interest to the accelerator scientists. The fact that the second beam shows higher or equal acceleration at the last timestep of the simulation is due to the fact that the first beam will outrun the wave later in time and therefore switches into a phase of deceleration while the second beam is still in an acceleration phase. In practice, when researchers want to select only the first beam, they usually perform selection of the particles using thresholding in px at an earlier time (e.g. $t = 27$) when the beam-particles of interest have maximum momentum, rather than the last timestep, here $t = 37$. By performing selection at an earlier time, one avoids selecting particles in the second beam while being sure to select all particles in the first beam. In this specific use case we are interested in analyzing and comparing the evolution of these two beams, which is why we performed selection at the last timestep.

C. Beam Formation

Having defined the beam in the 2D dataset, we can analyze formation of the beam by tracing the selected particles back to the time when particles entered the simulation and were then injected into the beam. In Figure 6, the individual particles of the beam are shown at timesteps $t = 14$ to $t = 17$. Here, color indicates the particles' momentum in the x direction (px). Figure 7 shows the particle traces over time colored according to particle ID's (the first beam appearing in blue and the second beam in yellow/red). Different sets of injection are readily visible, and two sets of particles appear at $t = 14$. The left bunch will be injected to form a beam in the second wake period, which is visible at $t > 14$. A second group of particles is just entering the right side of the box (recall that the simulation box is sweeping from left to right with the laser, so that plasma particles enter it from its right side). This bunch continues to enter at $t = 15$, and the particles stream into the first (rightmost) wake period. These particles are accelerated and appear as a bunch in the first bucket for $t > 15$. At the following timesteps $t = 16$ and $t = 17$, further acceleration of the two particle beams can be seen while only a few additional

particles are injected into the beam, and these are less focused, i.e., they show higher spread in the transverse direction (y).

D. Beam Refinement

Based on the information at timestep $t = 14$, we then refine our initial selection of the beam. By applying an additional threshold in x , we can select those particles of the beam that are injected into the first wake period behind the laser pulse (see Figure 8a and 8b). By comparing the temporal traces of the selected particle subset (green) with the traces of the whole beam (red), we can readily identify important characteristics of the beam (see Figure 8c). After being injected, the selected particle subset (green) first defines the outer part of the first beam at timestep $t = 15$, while additional particles are injected into the center of the beam. Later on at timesteps $t = 16$ and $t = 17$, the selected particles become strongly focused and define the center of the first beam. By refining selections based on information at an earlier time, we are able to identify characteristic substructures of the beam.

E. Beam Evolution

Using temporal parallel coordinates, we can analyze the general evolution of the beam in multiple dimensions (see Figure 9). Along the x axis, two separate beams can be seen at all timesteps ($t = 14$ to $t = 22$) with a quite stable relative position in x (x_{rel}). At early timesteps, both beams show similar acceleration in px while later on, at timestep $t = 18$ to $t = 22$ (magenta to lilac), the particles of the first beam show significantly higher acceleration with a relatively low energy spread. We find particles at relatively high relative positions in x direction (x_{rel}) only at timesteps $t = 14$ and $t = 15$ due to the fact that the particles of the beams enter the simulation window at these times.

F. 3D Analysis Example

We now describe a similar example analysis of the 3D particle dataset. Figure 10(a and b) shows the beam selection step for this dataset. At a much earlier timestep $t = 12$ ($x \approx 5.7 \times 10^{-4}$ compared to $x \approx 1.3 \times 10^{-3}$ in the 2D case) particles

are trapped and accelerated. In order to get an overview of the main relevant data, the user removed the background from the data first by applying a threshold of $px > 2.0 \times 10^9$. The user then selected particles in the first bunch via thresholding based on the momentum in x direction ($px > 4.586 \times 10^{10}$) and x position ($x > 5.649 \times 10^{-4}$) to exclude particles in the secondary periods from the selection. Figure 10b shows a volume rendering of the plasma density along with the selected particles revealing the physical location of the selected beam within the wake.

Figure 10c shows the traces of the particles selected earlier in Figure 10a. We selected the particles at timestep $t = 12$ then traced them back to timestep $t = 9$ where most of the selected particles enter the simulation window and forward in time to timestep $t = 14$. As one can see in the plot, the selected particles are constantly accelerated over time.

V. PERFORMANCE EVALUATION

In this section, we present results of a study aimed at characterizing the performance of our implementation under varying conditions using standalone, benchmark applications. These unit tests reflect the different stages of processing we presented earlier in Section IV, the Use Case study. First, we examine the serial performance of histogram computation in Section V-A: histograms serve as the basis for visually presenting data to a user via a parallel coordinates plot. Second, in Section V-B, we examine the serial performance of particle selection based on ID across all time steps of simulation data. Finally in Section V-C, we examine the parallel scalability characteristics of our histogram computation and particle tracking implementations on a Cray XT4 system.

For the serial performance tests in Sections V-A and V-B, we use the 3D dataset described earlier in Section IV: 30 timesteps worth of accelerator simulation data, each timestep having about 90 million particles and being $\approx 7GB$ in size (including $\approx 2GB$ for the index). The aggregate dataset size is about 210GB, including the index data. We store and retrieve simulation and index data using HDF5 and a veneer library called *HDF5-FastQuery* [32]. HDF5-FastQuery presents an implementation-neutral API for performing queries and obtaining histograms of data. We conduct the serial performance tests on a workstation equipped with a 2.2GHz AMD Opteron CPU, 4GB of RAM and running the SuSE Linux distribution.

The serial performance tests in Sections V-A and V-B measure the execution time of two different implementations that are standalone applications we created for the purpose of this performance experiment. One application, labeled “FastBit” in the charts, uses FastBit for index/query and histogram computation. The other, labeled “Custom” in the charts, does not use any indexing structure, and therefore performs a sequential scan of the dataset when computing histograms and particle selections. Note, in order to enable a fair comparison we here use our own Custom code which shows better performance than the IDL scripts currently used by our science collaborators.

A. Computing 2D Histograms

We run a pair of test batteries aimed at differentiating performance characteristics of computing both unconditional and conditional histograms. The unconditional histogram is simply a histogram of an entire dataset using a set of application-defined bin boundaries. A conditional histogram is one computed from a subset of data records that match an external condition.

1) *Unconditional Histograms*: In our use model, the computation of the unconditional histogram is a “one-time” operation. It provides the initial context view of a dataset.

For this test, we vary the number of bins in a 2D histogram over the following bin resolutions: 32×32 , 64×64 , 128×128 , 256×256 , 512×512 , 1024×1024 , and 2048×2048 bins. All histograms span the same range of data values: increasing the bin count results in bins of finer resolution. The results of this test are shown in Figure 11. Since both the FastBit and Custom applications need to examine all the data records to compute an unconditional histogram, we do not expect much performance variation as we change the number of bins. FastBit is generally faster than the custom code throughout primarily because of the difference in organization of the histogram bin counts array. FastBit uses a single array, which results in a more favorable memory access pattern. FastBit computes adaptive histograms by first computing a higher-resolution uniformly binned histogram and then merging bins. Since merging bins is a fairly inexpensive process and increasing the number of bins has no significant effect on the performance of uniformly binned histograms we observe only a minor, constant increase in computation time for adaptive versus uniform binning.

2) *Conditional Histograms*: In contrast to the unconditional histogram, which is a one-time computation, the process of visual data exploration and mining relies on repeated computations of conditional histograms. Therefore, we are particularly interested in achieving good performance of this operation to support interactive visual data analysis. This set of tests focuses on a use model, a change in the set of conditions results in examining a greater or lesser number of data records to compute a conditional histogram. This set of conditions reflects the repeated refinement associated with interactive, visual data analysis.

We parameterize our test based on the number of hits resulting from a range of different queries. We use thresholds of the form $px > \dots$ as the histogram conditions. As we increase the threshold of particle momentum in the x direction with this condition, fewer data records (particles) satisfy that condition and contribute to the resulting histogram. In these tests, we hold the number of bins constant at 1024×1024 .

Figure 12 presents results for computing conditional histograms using FastBit and the Custom application. We observe that for small number of hits, the FastBit execution times are dramatically faster than the Custom code that examines all data records. Also, in this regime we observe that uniform and adaptive binning show similar performance. While for unconditional histograms the minimum and maximum values are known for each variable, we here need to compute these

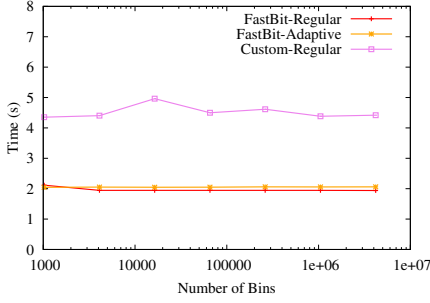


Fig. 11. Timings for serial computation of unconditional histograms.

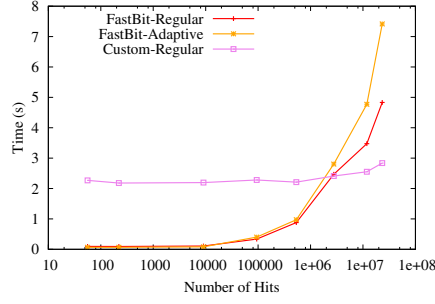


Fig. 12. Timings for serial computation of conditional histograms.

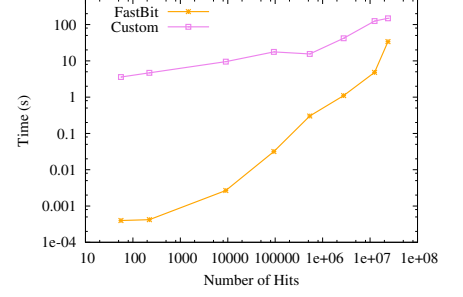


Fig. 13. Timings for serial processing of identifier queries.

values from the selected data parts in order to compute the adaptive binning. Due to this fact we observe a performance-decrease of the adaptive binning compared to uniform binning for very large selections. It is important to note that visual analysis queries typically isolate a small number of particles – from tens to thousands – and FastBit provides outstanding performance in this regime. As the number of hits increases, approaching 1M and 10M+ particles – which is a significant fraction of the 90M total particles – a sequential scan through all the data records produces better results. This performance change is due to the fact that FastBit computes the histogram in two separate steps. It first evaluates the user-specified conditions to select the appropriate values, and then counts the number of values in each histogram-bin. The selected values are passed from the first to the second step as an intermediate array with as many elements as the number of hits. It is expensive to pass this intermediate array through memory when it is large. Since the intended applications primarily have a small number of hits, using FastBit is more efficient.

B. Particle Selections

Following our use model, once a user has determined a set of interesting data conditions, like particles having a momentum exceeding a given threshold, the next activity is to extract those particles from the large dataset for subsequent analysis. This set of tests aims to show performance of the particle subsetting part of the processing pipeline.

The execution time for this task is clearly proportional to the size of the selection – the time required to find a set of particles in a large, time-varying dataset varies as a function of the size of the particle search set as well as the size of the simulation data itself. We parameterize the search set size for this set of performance tests, varying the number of particles to search for over values ranging from 10, 100, 1000, ... , up to 20M particles.

Our custom code uses a sequential scan of the entire dataset to search for particles in the search set. For each data record, it compares the particle ID of the record to the search set using an efficient algorithm: if the size of the search set is S , then the search time is $O(\log(S))$. If there are N data records in the entire dataset, the computational complexity of the entire algorithm is $O(N\log(S))$. In contrast, the worst-case

time required to locate a set of identifiers using FastBit is expected to be proportional to the number records found [17].

Figure 13 presents results for running ID queries using FastBit and the Custom code for one timestep. For relatively small numbers of identifiers, we observe that FastBit is about four orders of magnitude ($10^4\times$) faster than the Custom code. As the number of identifiers involved increases, the relative difference becomes smaller. When 20 million of identifies are involved, FastBit is still three times faster.

C. Scalability Tests

Whereas the previous sections have focused on the serial performance of computing histograms and performing particle subset selections, this section focuses on the scalability characteristics of both algorithms.

The platform for these tests is `franklin.nersc.gov`, a 9,660 node, 19K core Cray XT4 system. Each of the nodes consists of a 2.6GHz, dual-core AMD Opteron processor and has 4GB of memory and runs the Compute Node Linux distribution. One optimization we use on this machine at all levels of parallelism is to restrict operations to a single core of each node. This optimization maximizes the amount of memory and I/O bandwidth available to each process in our parallel performance tests. On this platform, the Lustre Parallel Filesystem serves data to each of the nodes. The nodes used in the study are a small fraction of a larger shared facility with a dynamic workload. Our scalability tests cover parallelism levels over the following range: 1, 2, 5, 10, 20, 50, and 100 nodes of the Cray XT4.

As in previous sections, we report wall-clock times which encapsulate CPU processing and I/O. We report the speedup factor as the ratio of time taken by a single node to the time taken by the node subset to complete a task. We do a strong scaling test by keeping the problem size fixed at 100 timesteps for all cases.

The dataset used in this study has 100 timesteps; each timestep has 177 million particles and is about 10GB in size. The aggregate dataset size is about 1.5TB, including the index data. We employ a fairly simple form of data partitioning for these scalability tests: namely we assign data subsets corresponding to individual timesteps (corresponding to individual HDF5 files) to individual nodes for processing. The subsets are statically assigned to nodes in a strided fashion.

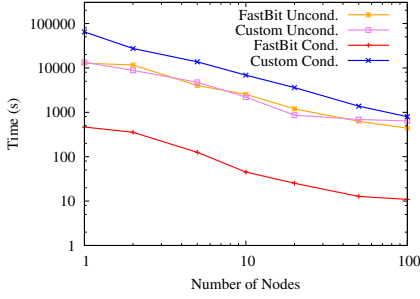


Fig. 14. Timings for parallel computation of histograms.

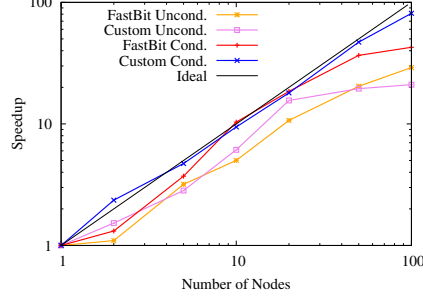


Fig. 15. Scalability of parallel histogram computation.

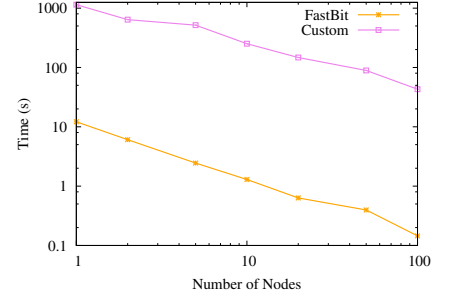


Fig. 16. Timings for parallel particle tracking.

Following the theme of use cases reported in previous sections, we report times for realistic science usage scenarios. For histogram computation, we generate five parallel histogram pairs for the position and momentum fields. We use 1024×1024 bins, which is a reasonable upper limit given typical screen resolutions. For conditional histograms, we use a $px > 7 \times 10^{10}$ query; for particle tracking, we report numbers for a $px > 10^{11}$ query, which results in 500 hits. All of these choices are grounded in discussions with our science collaborators and reflect reasonable ranges and thresholds.

Figure 14 presents results from parallelizing the computation of both conditional and unconditional histograms over multiple nodes. As expected, we observe that the computation time decreases as we add more nodes. Similar to the serial case, we do not observe much of a difference between FastBit and the Custom application for computing the unconditional histogram since both implementations examine all data records. For computing the conditional histogram, FastBit maintains its advantage over the Custom application. Figure 15 presents the speedup factors corresponding to this computation; we observe a very favorable speedup. This is to be expected since the nodes can perform their computations independently of others.

Figure 16 presents results from particle tracking over 100 timesteps. Similar to the serial case, FastBit is faster than the Custom application, and maintains its advantage when executed in parallel. Figure 17 demonstrates that we achieve excellent scalability in parallelizing this computation. We observe that when using 100 nodes, FastBit is able to track 500 particles over 1.5TB of data in 0.15 seconds. In contrast, the IDL scripts currently used by our collaborators take ≈ 2.5 hours to track 250 particles on a small 5GB dataset. Owing to IDL limitations, the scripts simply fail to run on the larger datasets used in this study. Our research marks a dramatic improvement for the scientists' workflow: what formerly required many hours can now be processed in less than a second.

VI. CONCLUSIONS AND FUTURE WORK

Managing and gaining insight from vast amounts of data is widely accepted as one of the primary bottlenecks in modern science. The work we present takes aim at enabling rapid knowledge discovery from large, complex, multivariate

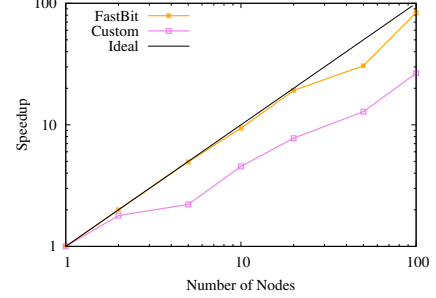


Fig. 17. Scalability of parallel particle tracking.

and time-varying scientific datasets and using modern HPC platforms. To achieve this objective, we have presented a novel approach for quickly creating histogram-based parallel coordinates displays. Further, we leverage this form of visual information display as the basis for forming complex multivariate range queries. We combined this form of visual information display with state-of-the-art index/query technology to quickly compute conditional histograms as well as to quickly and efficiently extract subsets of data that meet a set of conditions. We presented a case study showing how these technologies are used in concert to explore a large, time-varying dataset produced by a laser wakefield accelerator simulation. That exploration reveals interesting beam characteristics, such as how particles are first accelerated by the wakefield, then later “outrun” the wave and undergo a period of deceleration.

We also conducted a performance study of each of the fundamental algorithms that form a complete implementation in a production-quality, parallel capable visual data analysis application. That study shows the efficacy of our algorithms for computing histograms and for performing particle subset selection. These algorithms were shown to have favorable scalability characteristics over a set of processor pool sizes ranging from 1 up to 100 on a modern HPC platform, a Cray XT4. From the scientists' point of view, they are now able to perform a type of visual data analysis in a few seconds with this new work as compared to the several hours required using legacy tools. These results are significant because they show the successful synergy that can result when combining visual data analysis with scientific data management technologies to solve challenging problems in scientific knowledge discovery.

As we move forward with this work, we are particularly

interested in exploring different avenues for parallelizing the most expensive parts of the visual data analysis work. As visual data analysis is typically an interactive process, minimizing the response time is crucial to maximize its effectiveness. We also plan to apply this system to larger datasets from this and other scientific disciplines. As our implementation already supports a mode of operation whereby a parallel back end runs on the HPC platform to perform I/O, visualization and analysis processing with results being transmitted over the network to a client running on a workstation, we are in a good position to tackle such future challenges. Finally, the work we show here is primarily visual data analysis; we intend to extend this work to couple it with more traditional data analysis techniques.

ACKNOWLEDGMENT

This work was supported by the Director, Office of Science, Offices of High Energy Physics and Advanced Scientific Computing Research, of the U.S. Department of Energy under Contract No. DE-AC02-05CH11231 through the Scientific Discovery through Advanced Computing (SciDAC) program's Visualization and Analytics Center for Enabling Technologies (VACET) and the COMPASS project. This research used resources of the National Energy Research Scientific Computing Center (NERSC), which is supported by the Office of Science of the U.S. Department of Energy under Contract No. DE-AC02-05CH11231. We also thank the VORPAL development team for ongoing efforts to develop and maintain this parallel C++ framework on a variety of supercomputing platforms including those at NERSC.

REFERENCES

- [1] C. Nieter and J. R. Cary, "VORPAL: A Versatile Plasma Simulation Code," *J. Comput. Phys.*, vol. 196, no. 2, pp. 448–473, 2004.
- [2] C. Geddes, C. Toth, J. van Tilborg, E. Esarey, C. Schroeder, D. Bruhwiler, C. Nieter, J. Cary, and W. Leemans, "High-Quality Electron Beams from a Laser Wakefield Accelerator Using Plasma-Channel Guiding," *Nature*, vol. 438, pp. 538–541, 2004, IBNL-55732.
- [3] A. Inselberg, "Parallel coordinates for multidimensional displays," in *Spatial Information Technologies for Remote Sensing Today and Tomorrow, The Ninth William T. Pecora Memorial Remote Sensing Symposium*. IEEE Computer Society Press, 1984, pp. 312–324.
- [4] E. J. Wegman, "Hyperdimensional data analysis using parallel coordinates," *Journal of the American Statistical Association*, vol. 85, no. 411, pp. 664–675, Sep. 1990.
- [5] A. Inselberg, H. Hauser, M. Ward, and L. Yang, "Modern parallel coordinates: from relational information to clear patterns, tutorial," in *IEEE Visualization*, October 2006.
- [6] A. Inselberg, *Parallel Coordinates Visual Multidimensional Geometry and Its Applications*. Springer-Verlag, 2008.
- [7] Y.-H. Fua, M. O. Ward, and E. A. Rundensteiner, "Hierarchical parallel coordinates for exploration of large datasets," in *IEEE Visualization 1999*. Los Alamitos, CA, USA: IEEE Computer Society Press, 1999, pp. 43–50.
- [8] J. Johansson, P. Ljung, M. Jern, and M. Cooper, "Revealing structure within clustered parallel coordinates displays," in *INFOVIS '05: Proceedings of the Proceedings of the 2005 IEEE Symposium on Information Visualization*. Washington, DC, USA: IEEE Computer Society, 2005, p. 17.
- [9] M. Novotný, "Visually effective information visualization of large data," in *Proceedings of Central European Seminar on Computer Graphics (CESCG)*, 2004.
- [10] M. Novotný and H. Hauser, "Outlier-preserving focus+context visualization in parallel coordinates," *IEEE Transactions on Visualization and Computer Graphics*, vol. 12, no. 5, pp. 893–900, 2006.
- [11] S. Chaudhuri and U. Dayal, "An overview of data warehousing and OLAP technology," *ACM SIGMOD Record*, vol. 26, no. 1, pp. 65–74, Mar. 1997.
- [12] P. O'Neil, "Model 204 architecture and performance," in *2nd International Workshop in High Performance Transaction Systems, Asilomar, CA*, ser. Lecture Notes in Computer Science, vol. 359. Springer-Verlag, Sep. 1987, pp. 40–59.
- [13] C.-Y. Chan and Y. E. Ioannidis, "Bitmap index design and evaluation," in *SIGMOD*, 1998, pp. 355–366.
- [14] K. Wu, E. Otoo, and A. Shoshani, "On the performance of bitmap indices for high cardinality attributes," in *VLDB*, 2004, pp. 24–35.
- [15] FastBit is available from <https://codeforge.lbl.gov/projects/fastbit/>.
- [16] K. Wu, E. Otoo, and A. Shoshani, "Compressing bitmap indexes for faster search operations," in *SSDBM'02*, Edinburgh, Scotland, 2002, pp. 99–108.
- [17] —, "Optimizing bitmap indices with efficient compression," *ACM Transactions on Database Systems*, vol. 31, pp. 1–38, 2006.
- [18] E. W. Bethel, S. Campbell, E. Dart, K. Stockinger, and K. Wu, "Accelerating Network Traffic Analysis Using Query-Driven Visualization," in *Proceedings of 2006 IEEE Symposium on Visual Analytics Science and Technology*. IEEE Computer Society Press, October 2006, pp. 115–122, IBNL-59891.
- [19] K. Wu, W. Kogler, J. Chen, and A. Shoshani, "Using bitmap index for interactive exploration of large datasets," in *SSDBM'2003*. Washington, DC, USA: IEEE Computer Society, 2003, pp. 65–74.
- [20] K. Stockinger, E. W. Bethel, S. Campbell, E. Dart, and K. Wu, "Detecting Distributed Scans Using High-Performance Query-Driven Visualization," in *SC '06: Proceedings of the 2006 ACM/IEEE Conference on High Performance Computing, Networking, Storage and Analysis*. IEEE Computer Society Press, October 2006.
- [21] K. Stockinger, K. Wu, and A. Shoshani, "Strategies for processing ad hoc queries on large data warehouses," in *DOLAP'02*, McLean, Virginia, USA, 2002, pp. 72–79.
- [22] K. Wu, K. Stockinger, and A. Shoshani, "Breaking the curse of cardinality on bitmap indexes," in *SSDBM 2008*, 2008, pp. 348–365.
- [23] K. Stockinger, J. Shalf, K. Wu, and E. W. Bethel, "Query-Driven Visualization of Large Data Sets," in *Proceedings of IEEE Visualization 2005*. IEEE Computer Society Press, October 2005, pp. 167–174, IBNL-57511.
- [24] R. Bellman, *Adaptive Control Processes: A Guided Tour*. Princeton University Press, 1961.
- [25] K. Stockinger, E. W. Bethel, S. Campbell, E. Dart, and K. Wu, "Detecting Distributed Scans Using High-Performance Query-Driven Visualization," in *SC '06: Proceedings of the 2006 ACM/IEEE Conference on High Performance Computing, Networking, Storage and Analysis*. New York, NY, USA: IEEE Computer Society Press, October 2006, IBNL-60053.
- [26] 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*, October 2005, pp. 190–198.
- [27] VisIt is available from <https://wci.llnl.gov/codes/visit/>.
- [28] C. C. Law, A. Henderson, and J. Ahrens, "An application architecture for large data visualization: a case study," in *PVG 01: Proceedings of the IEEE 2001 symposium on parallel and large-data visualization and graphics*. IEEE Press, 2001, pp. 125–128.
- [29] EnSight Gold: <http://www.ensight.com/ensight-gold.html>.
- [30] S. Guha, K. Shim, and J. Woo, "Rehist: relative error histogram construction algorithms," in *VLDB '04: Proceedings of the Thirtieth international conference on Very large Data Bases*. VLDB Endowment, 2004, pp. 300–311.
- [31] C. G. R. Geddes, "Plasma channel guided laser wakefield accelerator," Ph.D. dissertation, University of California, Berkeley, 2005.
- [32] L. Gosink, J. Shalf, K. Stockinger, K. Wu, and E. W. Bethel, "HDF5-FastQuery: Accelerating Complex Queries on HDF Datasets using Fast Bitmap Indices," in *Proceedings of the 18th International Conference on Scientific and Statistical Database Management*. IEEE Computer Society Press, July 2006, IBNL-59602.