

Lawrence Berkeley National Laboratory

LBL Publications

Title

Feature Analysis, Tracking, and Data Reduction: An Application to Multiphase Reactor Simulation MFiX-Exa for In-Situ Use Case

Permalink

<https://escholarship.org/uc/item/0d30n31s>

Journal

Computing in Science & Engineering, 23(1)

ISSN

1521-9615

Authors

Biswas, Ayan
Ahrens, James P
Dutta, Soumya
et al.

Publication Date

2021

DOI

10.1109/mcse.2020.3016927

Peer reviewed

Feature Analysis, Tracking, and Data Reduction: An Application to Multiphase Reactor Simulation MFiX-Exa for In-Situ Use Case

Ayan Biswas

Los Alamos National Laboratory

Terece L. Turton

Los Alamos National Laboratory

James P. Ahrens

Los Alamos National Laboratory

Soumya Dutta

Los Alamos National Laboratory

Jordan M. Musser

National Energy Technology Laboratory

Ann S. Almgren

Lawrence Berkeley National Laboratory

Abstract—As we enter the exascale computing regime, powerful supercomputers continue to produce much higher amounts of data than what can be stored for offline data processing. To utilize such high compute capabilities on these machines, much of the data processing needs to happen in situ, when the full high-resolution data is available at the supercomputer memory. In this article, we discuss our MFiX-Exa simulation, which models multiphase flow by tracking a very large number of particles through the simulation domain. In one of the use cases, the carbon particles interact with air to produce carbon dioxide bubbles from the reactor. These bubbles are of primary interest to the domain experts for these simulations. For this particle-based simulation, we propose a streaming technique that can be deployed in situ to efficiently identify the bubbles, track them over time, and use them to down-sample the data with minimal loss in these features.

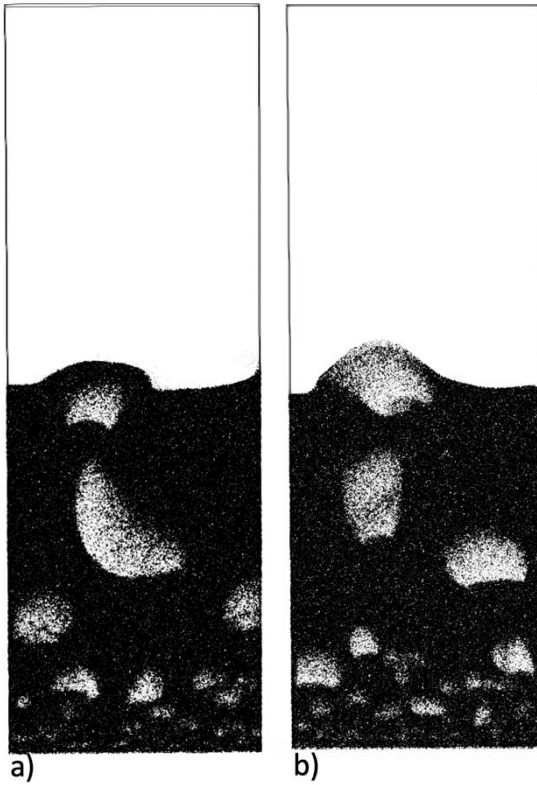


Figure 1. Example outputs from MFiX-Exa simulation. All the particles are shown for a) time step 150 and b) time step 408.

■ **EVERY YEAR**, supercomputers continue to become more powerful in terms of their computing capabilities. Compared to such high compute power, the I/O capabilities of the machines have not scaled proportionately. This creates a mismatch between the amount of data that can be computed by these supercomputers and how much data can be moved over to persistent storage for future data analysis and exploration tasks. Therefore, most of the important data analysis tasks for large-scale scientific simulations must occur while the data is still available at the supercomputer memory. Compared to the post-hoc analysis pipeline, this new streaming data analysis regime is referred to as *in-situ* or *online* processing.

For scientific simulations, data can be divided primarily into two groups: (1) grid-based and (2) particle-based. For grid-based data, generally the values of the variables (e.g., pressure, temperature) are defined on grid locations over a computational

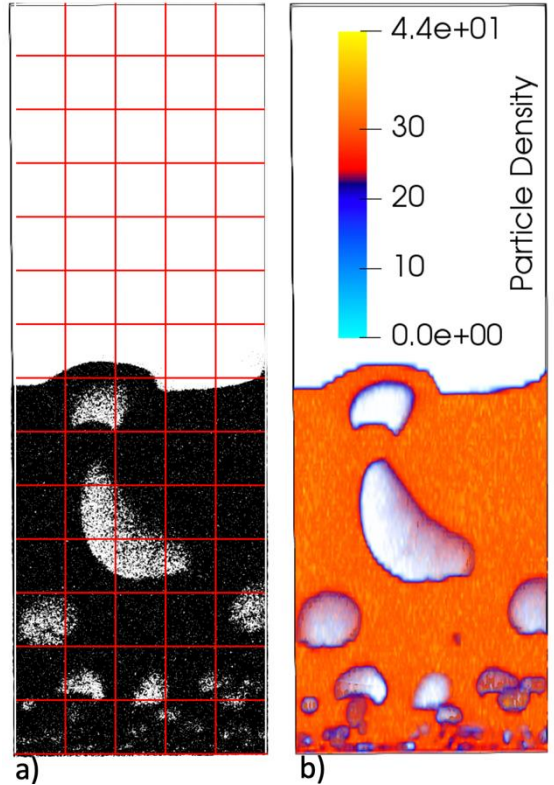


Figure 2. Conversion from particle dataset to density field. a) For timestep 150, an example of spatial histogram creation is shown. b) The visualization of the same data after conversion to its density field using a higher number of histogram bins. The high-density regions are shown in red.

domain. For particle-based data, a collection of particles is simulated by the physics codes that move through space and time. Depending on the simulation, individual particles with some local properties (e.g., energy transport via eddies in flow simulations) or a collection of particles with some global properties (e.g., particles forming halos in cosmological simulations) can be of interest to the domain experts.

In this article, we focus on a particle-based simulation, named MFiX-Exa. Starting from a set of particles, the domain experts for this simulation want to formulate and extract bubble statistics in each time step. Given that these bubbles get created, move over time, and get released into the air,

domain experts need to have efficient tracking of these features over time. But lots of particles are tracked over these time steps, and experts want to reduce the number of particles without losing the bubbles from the data. Due to the scale of the simulation, all these data analyses and reductions need to happen in batches, as the data becomes available in-situ. In this work, we propose a lightweight streaming bubble detection and tracking method and demonstrate a bubble-preserving particle-sampling algorithm. This algorithm is designed to handle the in-situ use case where data is distributed across multiple processors and the time steps are generated and analyzed in a streaming setting.

The manuscript is organized as follows: In the next section, we provide a brief overview of the existing works that relate to this paper. Next, we discuss the simulation and its data output. Then, we provide the details of our in-situ feature (bubbles, in this case) detection method. Next, we discuss our in-situ feature tracking algorithm. After that, we provide details on our feature-driven particle down-sampling approach. Finally, we conclude this manuscript with a summary of the proposed in-situ approaches.

RELATED WORK

The need for data reduction techniques for an in-situ environment has grown significantly in recent years, and various approaches to in-situ data reduction techniques have been proposed. While several scientists have explored data sub-sampling as one of the potential techniques [7], others have used distribution-based data representations as a suitable approach for in-situ data reduction [10]. Also, data reduction techniques such as wavelet-based data modeling techniques and data compression techniques [2,8] have been shown effective. For a comprehensive list of data reduction techniques, please refer to the state-of-the-art report by Li et al. [3].

Feature tracking is an important data analysis problem for scientific datasets. One of the earliest works in feature tracking was by Samataney et al. [5] where feature correspondence was used to track the features over time. Volume overlapping-based feature tracking was proposed by Silver and Wang [6]. Feature tracking in particle datasets was

investigated by Sauer et al., and the tracking was done in joint particle/volume datasets [9]. A predictor–corrector based feature tracking algorithm was proposed by Muelder et al. [11], and an in-situ application of a similar method was demonstrated in [4]. A broad survey of feature tracking techniques can be found in [12] by Post et al. Compared to the above works, here we focus on feature tracking techniques that can be deployed in an in-situ setting.

DESCRIPTION OF THE DATASET

The simulation represents a small-scale, pseudo two-dimensional ($0.15\text{m} \times 0.0032\text{m} \times 0.0508\text{m}$) system where a constant density (1.205 kg/m^3), constant viscosity ($1.8 \times 10^{-5} \text{ Pa}\cdot\text{sec}$) gas is used to fluidize spherical particles of uniform size ($148 \times 10^{-6} \text{ m}$ diameter) and density (1300 kg/m^3). A pressure outflow is prescribed at the top of the domain, while a constant velocity gas inlet (0.0342 m/sec) is specified along the bottom. The remaining boundaries are modeled as solid walls. The domain is decomposed into $672 \times 10 \times 228$ computational cells and initialized with solids by randomly distributing approximately 3.6 million particles. As the simulation progresses, particles settle in the direction of gravity, forming a bubbling fluidized bed (**Figure 1**).

Time-averaged particle statistics (e.g., averaged particle velocities) and transient field properties (e.g., bubble size and velocity) are commonly used to compare experimental data and model results [1]. Similar analysis could be employed in assessing reactor designs whereby bubble (and/or cluster) statistics would provide engineers insight into flow phenomena that adversely impact multiphase reactor performance (e.g., gas by-passing via bubble formation).

BUBBLE DETECTION IN IN-SITU SCENARIO

For this dataset, the bubble features are vaguely defined as *low-density regions* or *voids* in the data. To detect such features, we first need to compute the density field from the point dataset. Since the application mandates the algorithms be suitable for an in-situ environment where the data is distributed across multiple processors, we propose a histogram-based density estimate.

Histograms are essentially *counts* of values falling into each bin, and these counts are efficiently

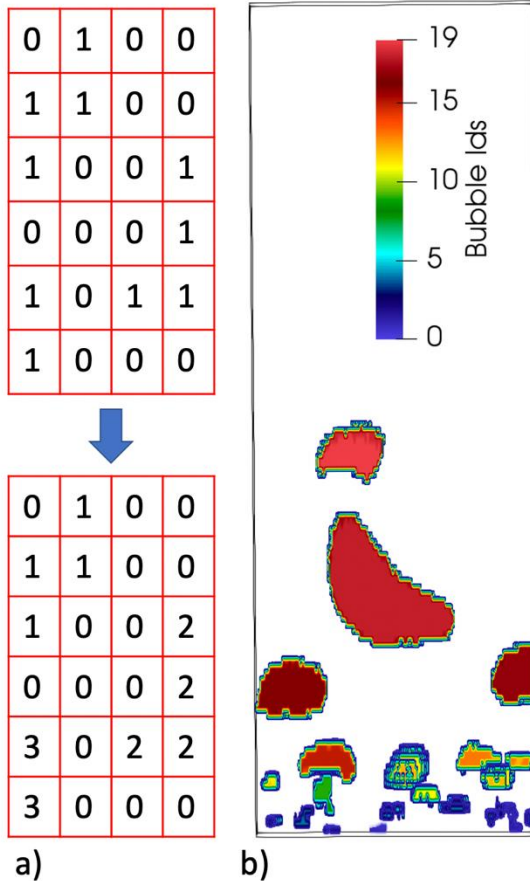


Figure 3. Individual bubble identification from the binary labeled field. a) Schematic example showing how a connected-components algorithm changes the binary field to assign individual bubbles a unique identifier at each time step. b) The visualization of time step 150 after identifying individual bubbles.

parallelizable because they can be added up across different processors. In our case, we use a histogram over the spatial domain. Each processor is responsible for the advection of a subset of particles, and it has the knowledge of the complete spatial domain of the simulation. For each processor, we divide the domain into the same three-dimensional regular grid (e.g., as shown in **Figure 2a**). Now, each processor can compute a local histogram by computing how many particles are falling into each bin. These local histograms can be summed together to derive the final density

histogram. Since the size of the histogram is very small compared to the size of the original data, we now perform all the operations on this histogram for bubble identification and tracking. Further, because each bin is spread across the domain, simply visualizing the density field will give us information regarding where the bubbles are located. We can think of this spatial histogram as a regular discretization of space, and it can be used as a regular grid dataset for analysis and visualization purposes.

For identification of the bubbles, we apply a thresholding (K) on the density field to produce our feature field ($F_{feature}$) which is essentially a binary version of the density field of same size:

$$F_{feature} = \begin{cases} 1 & \text{if } F_{density} \geq K \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

Essentially, for all locations where the density field value $F_{density} \geq K$, we set the locations to zero. Otherwise, the value is set to 1. This step assigns all the locations that are part of a bubble as 1, signifying that this is the feature of interest.

At this stage, although the bubbles can be visualized (**Figure 2b**), they cannot be distinguished from each other because all the bubbles are assigned a value of 1. To prepare each bubble for tracking and collecting statistics individually, we mark each bubble region with a unique identifier. To achieve this, we perform a *connected components* algorithm on this binary field (**Figure 3a**). Primarily used for images, this algorithm finds the total number of components in the field. We apply a flood-fill-based connected component algorithm for identifying all the bubbles. This method is performed on only one processor because we are not applying it to the particles, rather it is applied to the binary version of the density histogram.

After this stage, we need to make one correction to the detected bubbles. The top layer of air will be detected as one bubble, which in fact, should not be considered a bubble at all. To fix this, we determine which bubble has a top layer that is the same as the data domain. We remove that bubble from our list of bubbles (**Figure 3b**) and then move on to track them over time.

BUBBLE TRACKING IN IN-SITU SCENARIO

Understanding how bubbles evolve is an important aspect of the exploration process for this dataset. Bubbles get created; move over time; and merge, split, or die. Domain experts want to explore such events associated with the bubbles for the course of the simulation time. In the previous section, we described how the individual bubbles can be assigned a unique identifier. But when such a method is applied independently for each time step, the same bubble might be assigned a different identifier when it moves from time step t to $t + 1$. Tracking ensures that the same bubbles always get the same tag or id irrespective of which time step they belong to.

Tracking consists of the following scenarios:

1. Born: A new bubble is born.
2. Move: A bubble has moved to a new location with or without changing shape.
3. Split: An existing bubble has split into two or more smaller bubbles.
4. Merge: Two (or more) smaller bubbles have merged to form a larger bubble.
5. Die: A bubble moves out of the particles and gets mixed into the air.

Identification of these events in-situ can be challenging due to the streaming nature of the problem—the time steps come one at a time and as we move to the next time step, the previous time step no longer exists in the memory. To alleviate these constraints, we retain the bubble information from the previous time step to match with the next one. Because the bubble information generated from the spatial histograms is much smaller in size compared to the full-scale particle data, we can carry it from a given time step to the next.

To resolve the tracking scenarios (Figure 4a), we first try to detect the *move* event. For this event, we use the spatial overlap criteria. For a given bubble, we attempt to search for another bubble in the next time step that has maximum overlap. Because all the bubbles in this dataset are primarily moving upwards, we use the predictor-corrector method [11]. We first predict that the bubble will be moving upwards given the average speed of the surrounding particles of that location and then find the bubble with maximum spatial overlap with this predicted bubble location. When a *move* event is found, the bubble id from the previous time step is

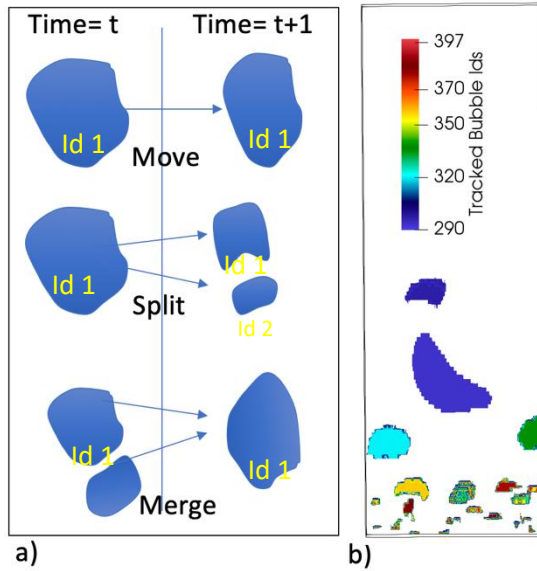


Figure 4.2 Illustration of proposed in-situ bubble tracking method. a) From one time-step to the next, existing bubbles can move (change shape), split, or merge with other bubbles. A new bubble can also be born with a new id. b) After applying bubble tracking, visualization from time step 150. Compared to Figure 3b, we see the bubble ids are now more numerous in Figure 4b because they are uniquely assigned to each bubble moving over time.

assigned to the *moved bubble* of the new time step. This way we can be consistent with the bubble ids over time.

Next, we check for split/merge events. For detecting split events, we check the bubble data of the current time step against the previous one. For each bubble of the current time step, we perform a location-based overlap comparison with the bubbles from the previous time step. If we observe there were overlaps with more than one bubble, we conclude a merge event occurred. Merge and split can be thought of as similar events except reversed in time direction; that is, a merge event in the forward comparison of time t with $t + 1$ can be thought of as a split event when comparing time $t + 1$ with time t . We can use this idea to detect merge events similar to detecting split events in situ, except we compare time t with time $t + 1$ for bubble overlaps, as mentioned in the split detection earlier.

After the application of this tracking algorithm, bubbles can be efficiently assigned corresponding

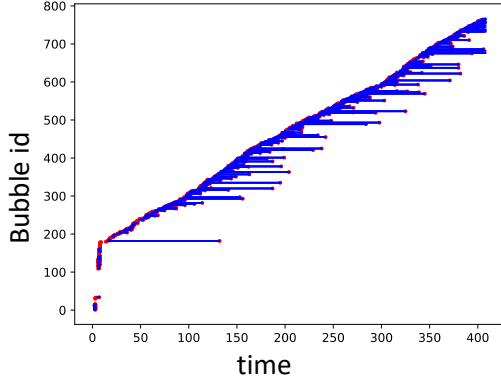


Figure 5. Bubble statistics: visualization of evolution of bubbles over time.

identifiers that are consistent over time. In **Figure 4b**, we show the outcome of applying this method on the MFiX-Exa time step 150. Compared to **Figure 3b**, where our method is not applied, we observe that the bubble ids are more numerous when using our tracking method. An increased number is intuitive given that after tracking bubbles for approximately 150 time steps, many bubbles have appeared or disappeared. Since unique ids are assigned over time steps, we see that the minimum bubble id in **Figure 4b** starts at 292, whereas in **Figure 3b**, all the bubbles have local ids, resulting in a range of numbers from zero to 19.

As we track the bubbles over time, we can further collect their statistics and life-span summary reports. For each bubble, we can track its time history (**Figure 5**) and compute properties—such as volume—as it moves through time.

BUBBLE-AWARE PARTICLE SAMPLING IN IN-SITU SCENARIO

For exascale simulations such as MFiX-Exa, one of the primary needs is to reduce the datasets such that they can also be explored in a post-hoc analysis pipeline. For in-situ data reduction, there are a few popular choices available to domain experts, such as random and regular sampling. For a particle dataset such as this one, regular sampling often generates less useful results because it can produce visualization artifacts. Random sampling, on the other hand, is generally a safe choice given that the selection of particles is random, thereby avoiding sampling artifacts.

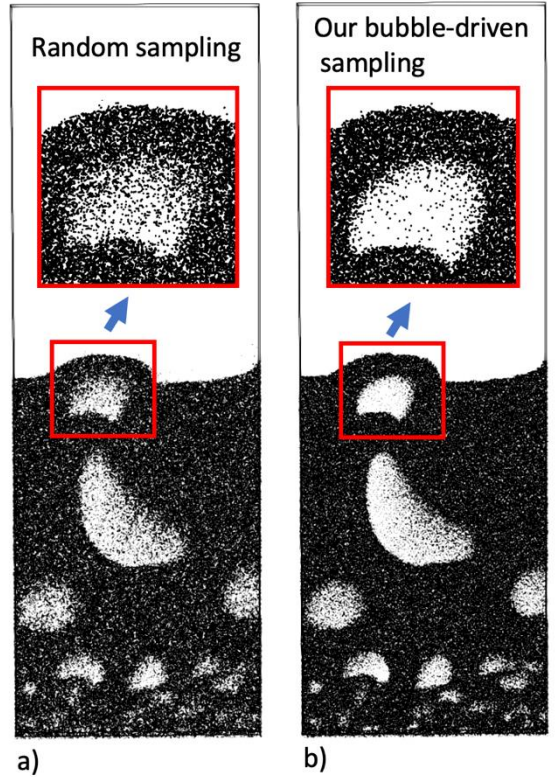


Figure 6. Sampling results for the particles from time step 150. Five percent of samples have been taken using a) random sampling and b) our proposed bubble-preserving sampling. As seen from the results, the bubbles are more prominent and better preserved (refer to the zoomed-in view) with our method.

One drawback of using such generic sampling methods for a dataset where the experts are interested in preserving the features (in this case, bubbles) stems from the fact that random methods cannot ensure the features of the dataset will be preserved with high fidelity. In this manuscript, we propose feature-based sampling for this dataset so that the features of the data are better preserved. In our approach, we first perform the bubble detection as mentioned in Section “BUBBLE DETECTION IN IN-SITU SCENARIO”. Since bubble information is contained in a histogram-like derived product, when performing sampling, we use this density histogram. To ensure the bubbles of the data are preserved as best as possible given the

storage limitation, our primary goal is to preserve the density histogram. For example, if the original time step consisted of 1 million particles and we need to select 100,000 particles from that time step, the selections are made such that the resulting density histogram from the sampled data looks very similar to the histogram from the original data. To ensure this, we devise a two-pass algorithm over the data where first the density histogram is created. Then, the bubble regions are left empty and particles from the non-bubble regions are selected proportionate to the original density. Again, since we are applying the analysis primarily to a histogram that is much smaller in size compared to the data, this histogram is shared across all the processors. This enables distributed sampling of the particles.

An example of such sampling is shown in **Figure 6**. In this case, we take 5% of samples from the original data using both the random method (Figure 6a) and our proposed method (Figure 6b). Looking closely into the zoomed-in views, it can be observed that the bubbles are better preserved with our data-sampling method.

DISCUSSION

In this section, we briefly discuss a few of the design choices and reproducibility aspects that will facilitate porting this workflow to the parallel and distributed setting.

Histogram Creation: Creation of density histogram from the particles is an important step. For scalability, histograms can be created efficiently by the local processors and then MPI_Reduce like operation can be applied to collect all the local histograms into a global histogram. This keeps the data communication quite low and scale to a large-scale data simulation.

Memory Footprint: Global histograms are used by each processor for sampling and bubble detection. Generally, given 3.6 million particles that need x,y,z locations to be stored, a density histogram bin resolution of 128x8x64 produces almost two orders of magnitude smaller memory overhead.

Data Storage: While sampling the particles, each processor operates individually, resulting in good scaling. When writing out the reduced set of particles, parallel I/O modules such as HDF5 can be incorporated in future.

Reproducibility: The software package MFiX is available at <https://mfix.netl.doe.gov> ([password-protected repository](#)); however, to date no public releases of MFiX-Exa have been made available. MFiX-Exa will be made publicly available in the future, but no release schedule has been established. The particle-tracking software was developed in Los Alamos National Laboratory and is currently in the process of being released as an open source code. Due to national laboratory rules, open sourcing of code is a time-consuming process. We would like to make the code available to the readers as soon as the open-source release process is completed.

FUTURE WORK

In the future, we plan to deploy these in-situ-ready algorithms in a real in-situ setting with MFiX-Exa simulations and conduct a detailed performance analysis. We would like to explore generic feature-detection approaches that can apply to a variety of particle/grid-based simulations. Additionally, we would like to apply both unsupervised (e.g., clustering) and semi-supervised (e.g., collecting some expert labeled data and training a classifier) machine learning techniques for more robust and automated feature extraction and tracking. Finally, we would like to explore the possibilities of new data reduction schemes that can be incorporated into the in-situ scenario.

CONCLUSION

In this manuscript, we discussed a feature-based analysis of MFiX-Exa particle simulations for various in-situ use cases. We described a lightweight bubble-detection algorithm based on spatial histograms for particle density computation. As the time steps arrive one by one in the in-situ scenario, we discussed a bubble-tracking approach from these individually detected bubbles. With a reliable bubble extraction method, we discussed how that can be used for collecting efficient samples from the original set of particles. This will enable in-situ data reduction with high fidelity feature preservation. Our proposed in-situ-ready approaches are primarily geared towards MFiX-Exa simulations, but it can still be useful to other particle-based simulations with similar feature-driven analysis needs.

ACKNOWLEDGMENT

The authors would like to thank the Department of Energy and Los Alamos National Laboratory for the funding and support in carrying out this research. This research was supported by the Exascale Computing Project (17-SC-20-SC), a collaborative effort of the U.S. Department of Energy Office of Science and the National Nuclear Security Administration.

■ REFERENCES

1. Müller, C. R., Scott, S. A., Holland, D. J., Clarke, B. C., Sederman, A. J., Dennis, J. S., & Gladden, L. F. (2009). Validation of a discrete element model using magnetic resonance measurements. *Particuology*, 7(4), 297-306.
2. P. Lindstrom, "Fixed-Rate Compressed Floating-Point Arrays," in *IEEE Transactions on Visualization and Computer Graphics*, vol. 20, no. 12, pp. 2674-2683, 31 Dec. 2014.
3. Li, S., Marsaglia, N., Garth, C., Woodring, J., Clyne, J. and Childs, H. (2018), Data Reduction Techniques for Simulation, Visualization and Data Analysis. *Computer Graphics Forum*, 37: 422-447.
4. P. N. Duque, D. E. Hiepler, S. M. Legensky and C. P. Stone, In-Situ Feature Tracking and Visualization of a Temporal Mixing Layer, 2012 SC Companion: High Performance Computing, Networking Storage and Analysis, Salt Lake City, UT, 2012, pp. 1593-1593.
5. R. Samtaney, D. Silver, N. Zabusky and J. Cao, "Visualizing features and tracking their evolution," in *Computer*, vol. 27, no. 7, pp. 20-27, July 1994.
6. D. Silver and X. Wang, "Volume tracking," *Proceedings of Seventh Annual IEEE Visualization '96*, San Francisco, CA, 1996, pp. 157-164.
7. Ayan Biswas, Soumya Dutta, Jesus Pulido, and James Ahrens. 2018. In situ data-driven adaptive sampling for large-scale simulation data summarization. In *Proceedings of the Workshop on In Situ Infrastructures for Enabling Extreme-Scale Analysis and Visualization (ISAV '18)*. Association for Computing Machinery, New York, NY, USA, 13–18.
8. W. Austin, G. Ballard and T. G. Kolda, "Parallel Tensor Compression for Large-Scale Scientific Data," *2016 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*, Chicago, IL, 2016, pp. 912-922, doi: 10.1109/IPDPS.2016.67.
9. F. Sauer, H. Yu and K. Ma, "Trajectory-Based Flow Feature Tracking in Joint Particle/Volume Datasets," in *IEEE Transactions on Visualization and Computer Graphics*, vol. 20, no. 12, pp. 2565-2574, 31 Dec. 2014.
10. S. Dutta, C. Chen, G. Heinlein, H. Shen and J. Chen, "In Situ Distribution Guided Analysis and Visualization of Transonic Jet Engine Simulations," in *IEEE Transactions on Visualization and Computer Graphics*, vol. 23, no. 1, pp. 811-820, Jan. 2017.
11. C. Muelder and K. Ma, "Interactive feature extraction and tracking by utilizing region coherency," *2009 IEEE Pacific Visualization Symposium*, Beijing, 2009, pp. 17-24.
12. Post, F.H., Vrolijk, B., Hauser, H., Laramée, R.S. and Doleisch, H. (2003), *The State of the Art in Flow Visualisation: Feature Extraction and Tracking*. *Computer Graphics Forum*, 22: 775-792.

Ayan Biswas is a scientist in the Data Science group (CCS-7) at Los Alamos National Laboratory, Los Alamos, New Mexico, USA. His current research focuses on large-scale data reduction, in-situ data analysis, uncertainty quantification, statistical analysis, and high-dimensional data visualization. He received his Ph.D. in Data Visualization from The Ohio State University. Dr. Biswas is a member of the IEEE and the IEEE Computer Society. Contact him at ayan@lanl.gov.

James P. Ahrens is a senior scientist at Los Alamos National Laboratory. His research interests include large-scale data analysis and visualizations. He received his Ph.D. from University of Washington. Dr. Ahrens is a member of the IEEE and the IEEE Computer Society. Contact him at ahrens@lanl.gov.

Soumya Dutta is a staff scientist in the Data Science at Scale team at Los Alamos National Laboratory. He received his M.S. and Ph.D. degrees in Computer Science and Engineering from The Ohio State University in May 2017 and May 2018, respectively. His research interests include statistical data modeling, summarization, and analysis; in-situ data analysis, reduction, and feature exploration; HPC; uncertainty quantification analysis; and time-varying, multivariate data exploration. Contact him at sdutta@lanl.gov.

Ann S. Almgren is a senior scientist in the Computational Research Division of Lawrence Berkeley National Laboratory and the Group Lead of the Center for Computational Sciences and Engineering. Her primary research interest is in computational algorithms for solving PDEs in a variety of application areas. Her current projects include the development and implementation of new multiphysics algorithms in high-resolution adaptive mesh codes designed for the latest hybrid architectures. She is an SIAM Fellow and the Deputy Director of the ECP AMR Co-Design Center. Contact her at asalmgren@lbl.gov.

Jordan Musser is a physical research scientist at the National Energy Technology Laboratory. His research interests include the development and application of computational multiphase models. He received his Ph.D. from West Virginia University. Contact him at jordan.musser@netl.doe.gov.

Terece L. Turton is a staff scientist at Los Alamos National Laboratory. Her current research interests

include perceptual user evaluation and workflow analysis in scientific visualization. She received her Ph.D. from University of Michigan. Contact her at tturton@lanl.gov.