

Benchmarking the Graphulo Processing Framework

Timothy Weale[†], Vijay Gadepally^{‡§*}, Dylan Hutchison[°] and Jeremy Kepner^{‡§+}

[†]Department of Defense, [‡]MIT Lincoln Laboratory, [§]MIT Computer Science & AI Laboratory

⁺MIT Mathematics Department, [°]University of Washington

Abstract—Graph algorithms have wide applicability to a variety of domains and are often used on massive datasets. Recent standardization efforts such as the GraphBLAS specify a set of key computational kernels that hardware and software developers can adhere to. Graphulo is a processing framework that enables GraphBLAS kernels in the Apache Accumulo database. In our previous work, we have demonstrated a core Graphulo operation called *TableMult* that performs large-scale multiplication operations of database tables. In this article, we present the results of scaling the Graphulo engine to larger problems and scalability when a greater number of resources is used. Specifically, we present two experiments that demonstrate Graphulo scaling performance is linear with the number of available resources. The first experiment demonstrates cluster processing rates through Graphulo’s *TableMult* operator on two large graphs, scaled between 2^{17} and 2^{19} vertices. The second experiment uses *TableMult* to extract a random set of rows from a large graph (2^{19} nodes) to simulate a cued graph analytic. These benchmarking results are of relevance to Graphulo users who wish to apply Graphulo to their graph problems.

I. INTRODUCTION

Large-scale graph analytics are of interest in a variety of domains. For example in an intelligence, surveillance, and reconnaissance (ISR) application, graph analytics may be used for entity matching or track analysis. For social media applications, graph analytics may be used to determine communities or patterns of interest. Recently, an effort known as the GraphBLAS [1], [2] (<http://graphblas.org/>) has sought to standardize a set of computational kernels and demonstrate their applicability to a wide variety of graph analytics [3]. The connection between graph processing and linear algebra is well defined in [4], [5]. The GraphBLAS specification [5]–[7] currently consists of five core operations: sparse matrix multiplication, element-wise matrix multiplication, matrix addition, reference, and dereference, along with a small number of additional helper functions.

Given that the vast majority of data are stored in databases such as Apache Accumulo, there is also interest in developing processing frameworks, such as GraphX [8], Pregel [9], Gaffer [10], and GraphLab [11], optimized for graph operations on data in databases. A survey of popular graph database models

is provided in [12]. It is natural that a processing framework designed for databases makes use of the GraphBLAS kernels.

Graphulo [13] is a specialized graph processing framework built to work with Apache Accumulo and to conform to the GraphBLAS standard. Apache Accumulo is a NoSQL database designed for high performance ingest and scans [14], [15]. As a general-purpose warehousing solution, Accumulo has wide adoption in a variety of government and non-government settings.

As shown in Figure 1, Accumulo’s data model has a key with 5 components: row, column family, column qualifier, visibility, and timestamp. For simplicity, Graphulo provides directed graph structures using the row (source vertex), column qualifier (destination vertex), and the value (the number of connections between the two vertices). The visibility and timestamp attributes are not investigated in this paper.

Key					Value
Row ID	Column			Timestamp	
	Family	Qualifier	Visibility		

Fig. 1. Accumulo Key-Value Data Model

One of Accumulo’s advantages over other data warehousing solutions is its *iterator* framework for in-database computation. This capability provides improved performance by exploiting data locality, avoiding unnecessary network transfers, and moving the processing to the stored data. Additionally, by co-locating storage and computation, one can take full advantage of other aspects of Accumulo’s infrastructure, such as write-ahead logging, fault-tolerant execution, and parallel load balancing of data. As a first-order operator, iterators have full access to Accumulo’s data model, allowing for highly flexible, server-side algorithm development capabilities.

In the remainder of this paper, we provide a brief high-level review of Graphulo in Section II, including a description of the Accumulo iterators used and the process for tracking the amount of work done in each experiment. Section III presents the main SpGEMM scaling experiment using Graphulo, including experimental setup and results. We then describe a row extraction experiment in Section IV, providing scaling results for another real-world Graphulo use case. Finally, we discuss related work and areas for optimization in Section V before concluding in Section VI.

*Vijay Gadepally is the corresponding author and can be reached at vijay[at]ll.mit.edu

This material is based upon work supported by the Assistant Secretary of Defense for Research and Engineering under Air Force Contract No. FA8721-05-C-0002. Any opinions, findings and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the Assistant Secretary of Defense for Research and Engineering.

II. GRAPHULO

At a high level, the goal of the Graphulo project is to allow developers to focus on analytic development, rather than on writing custom MapReduce processes for graph data stored in Apache Accumulo. Thus, Graphulo explicitly implements the GraphBLAS specification on top of Apache Accumulo [3]. In our experiments, we focus on testing the performance of the Sparse Generalized Matrix Multiply (SpGEMM), the core kernel at the heart of GraphBLAS, because SpGEMM can be readily extended (via user-defined multiplication and addition functions) to express many GraphBLAS primitives. Additionally, SpGEMM is an essential part of a wide range of algorithms, including graph search, table joins, and many others [4].

Given matrices \mathbf{A} and \mathbf{B} as well as operations \oplus and \otimes for scalar addition and multiplication, the matrix product $\mathbf{C} = \mathbf{A} \otimes \mathbf{B}$, or for convenience, $\mathbf{C} = \mathbf{AB}$, defines entries of result matrix \mathbf{C} as

$$\mathbf{C}(i, j) = \bigoplus_k \mathbf{A}(i, k) \otimes \mathbf{B}(k, j) \quad (1)$$

We call the implementation of SpGEMM in Accumulo *TableMult*, short for multiplication of Accumulo tables. Accumulo tables have many similarities to sparse matrices, although a more precise mathematical definition is associative arrays [16]. For this work, we concentrate on distributed tables that may not fit in memory and use a streaming approach that leverages Accumulo's built-in distributed infrastructure. In Figure 2, we present the Graphulo *TableMult* iterator stack that implements Equation 1. *TableMult* operates on a table storing the transpose \mathbf{A}^T and a table storing \mathbf{B} and writes its output to a separate table \mathbf{C} .

TableMult's execution involves three Accumulo iterators: the *RemoteSourceIterator*, the *TwoTableIterator* and the *RemoteWriteIterator*. The *RemoteSourceIterator* scans the entries of table \mathbf{A}^T and provides them as inputs to the *TwoTableIterator*, which aligns and creates partial products. This action is the implementation of the $\mathbf{A}(i, k) \otimes \mathbf{B}(k, j)$ portion of Equation 1. We call intermediary results of operations *partial products* (*pp*). For the sake of sparse matrices, we only perform \oplus and \otimes when both operands are nonzero. The processing rate of the Graphulo cluster is given by the number of partial products divided by the processing time.

The results of the multiplication are then written to table \mathbf{C} with the *RemoteWriteIterator*. The summation (\oplus operator) is accomplished through combiner iterator subclasses placed on \mathbf{C} . This lazy summation only runs when table \mathbf{C} is required to execute a scan or minor/major compaction and helps maximize throughput while guaranteeing correctness of output.

III. EXPERIMENT: GRAPH MULTIPLICATION

Prior work [13] compared Graphulo's performance against the Dynamic Distributed Dimensional Data model (D4M), demonstrating improvement in terms of both speed and scale of computation. In this experiment, we demonstrate Graphulo's ability to scale its processing with the number of

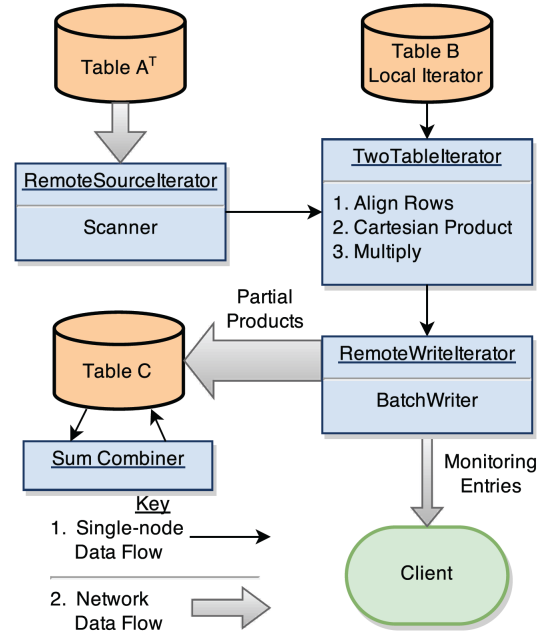


Fig. 2. Graphulo *TableMult* Iterator Processing Steps and Data Flow

cluster nodes. As in the previous experiment, we create and store random input graphs in our Accumulo tables using the Graph500 unpermuted power-law graph generator [17]. Power-law graphs have a highly connected initial vertex and exponentially decreasing degrees for subsequent vertices and have been used to describe a diverse set of real graphs [18]. The power-law generator takes *SCALE* and *EdgesPerVertex* as parameters, creating graphs with 2^{SCALE} rows and $EdgesPerVertex \times 2^{SCALE}$ entries. For our experiments, we fix *EdgesPerVertex* to 16 and use the *SCALE* parameter to vary problem size.

Our experimental setup is as follows:

- 1) Generate two graphs with different random seeds and insert them into Accumulo as adjacency tables.
- 2) For each table, identify optimal split points for each input graph. These optimal split points depend on the number of tablets,* which scales with the number of nodes in the cluster.
- 3) Set the input graphs' table splits equal to that point.
- 4) Create an empty output table. Pre-split it with an input split position from one of the two input tables.
- 5) Compact the input and output tables so that Accumulo redistributes the tables' entries into the assigned tablets.
- 6) Run and time Graphulo *TableMult* multiplying the transpose of the first input table with the second.
- 7) Repeat the experiment several times to derive a representative processing rate and minimize noise.

Previous work included an additional processing step to determine the optimal tablet split for the output graph. We

*In our configuration, a single-node cluster had each table split across two tablets; two nodes had four tablets; etc.

TABLE I
NUMERICAL RESULTS AND PARAMETERS FOR LARGE-GRAPH *TableMult* MULTIPLICATION EXPERIMENTS

Nodes	SCALE	Partial Products (pp)	Rate (pp/sec)	Speedup Ratio	SCALE	Partial Products (pp)	Rate (pp/sec)	Speedup Ratio
1	18x18	2.94×10^9	1.92×10^5	1.00	17x17	1.09×10^9	1.83×10^5	1.00
2			3.86×10^5	2.01	18x17	1.68×10^9	4.01×10^5	2.19
4			7.53×10^5	3.91	18x18	2.94×10^9	7.61×10^5	4.17
8			1.38×10^6	7.17	19x18	4.53×10^9	1.36×10^6	7.47
16			1.87×10^6	9.74	19x19	7.85×10^9	1.90×10^6	10.40

include no such “oracle” capability and set the output node distribution on the basis of one of the two input graphs, which can be calculated at ingest. As the number of nodes and tablets increased, calculating the optimal split became an overwhelming part of the experimental runtime.

The experiments were conducted on a shared system, with each node initially allocated 4GB of memory to an Accumulo tablet server (allowing growth in 1G steps), 1GB for Java-managed memory maps and 128MB for each data and index cache. The installed software was Graphulo 1.0.0, Accumulo 1.6.0, Hadoop 2.2.0 and ZooKeeper 3.4.6. Each table was split across two tablets per node, allowing for some additional processing parallelization.

Figure 3 shows the results for the first experiment. In this, the single-node cluster processing rate is used as a baseline for comparison. All other plots demonstrate the relative speedup from the single-node example. Ideal linear speedup is a line corresponding to a plot of $y = x$.

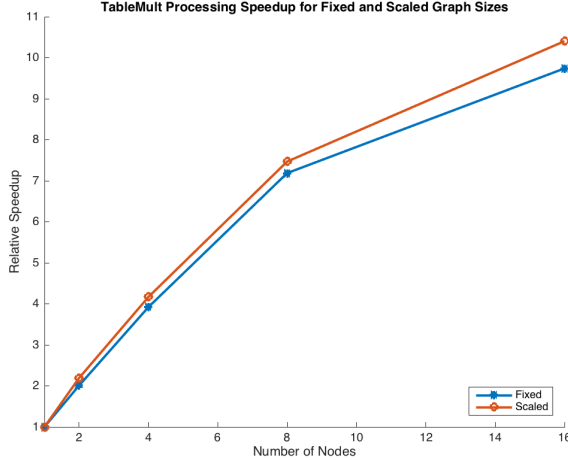


Fig. 3. Relative Processing Speedup for TableMult Operation Across Fixed-Size and Scaled Graphs

In the “Fixed” experiment, the size of the graphs involved in the TableMult operation were held constant at SCALE = 18. This strong-scale experiment allowed us to examine the processing speedup across the entire cluster by standardizing the total amount of work to be done and normalizing the entire cluster performance against a single node. The “Scaled” experiment increases the size of one of the graphs in the TableMult operation with the number of nodes. Increasing both

the graph size and number of nodes by a power of 2 allows for a roughly linear growth pattern that keeps the amount of work done per node roughly constant, providing weak scaling results for Graphulo.

The experiment demonstrates that Graphulo retains Accumulo’s ability to grow linearly with the number of processing nodes. In both the strong and weak scaling experiments, Graphulo closely mirrored the expected linear speedup. Numerical results of the first experiment can be found in Table I. Some tweaks (discussed below) should improve performance in the future.

IV. EXPERIMENT: SUBGRAPH EXTRACTION

Organizations often have a massive amount of information available in their databases. At any time, however, only certain subsets of interest will be required for analysis and processing. Graphulo can provide a quick way to extract and use these subsets, improving an algorithm’s ability to run effectively. We refer to such analytics as cued analytics, that is, analytics on selected table subsets. Such analytics are ideally suited for databases by quickly accessing subsets of interest (for example, particular communities in a large social media graph) in contrast to whole-table analytics that may be better suited for parallel processing frameworks such as Hadoop Map-Reduce [19].

In this experiment, we take a large sample graph (SCALE = 19) and randomly extract some portion for analysis (in our experiments, either 1024 or 2048 graph edges) by randomly generating a set of vertices to sample ($\{sampleset\}$). This random set of vertices is then used to create a binary diagonal matrix (using Equation 2) that contains 1s at the randomly selected vertex locations (Equation 3). These operations are illustrated below:

$$E = \begin{pmatrix} a_{1,1} & 0 & \cdots & 0 \\ 0 & a_{2,2} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & a_{n,n} \end{pmatrix} \quad (2)$$

$$a_{n,n} = \begin{cases} 1 & \text{if } n \in \{sampleset\} \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

The generated binary diagonal matrix (E) can then be used to extract the rows (G_{samp}) from the original graph (G_{orig}) by using a TableMult operation with the two graphs (E, G_{orig}). The experimental steps are as follows:

- 1) Insert a single, large graph into Accumulo as adjacency tables. This will be the graph for the entire experiment.
- 2) Generate a random diagonal matrix with the appropriate number of non zero entries (1024, 2048).
- 3) Extract a set of random rows using the Graphulo Table-Mult operation.
- 4) Repeat the random extraction process several times to derive a representative processing rate.

We use the same cluster configuration as in the prior experiment: a shared system with each node initially allocated 4GB of memory to an Accumulo tablet server (allowing growth in 1G steps), 1GB for Java-managed memory maps and 128MB for each data and index cache. The installed software was Graphulo 1.0.0, Accumulo 1.6.0, Hadoop 2.2.0 and ZooKeeper 3.4.6. Each table is split across 2 tablets per node to allow for additional processing parallelization.

As mentioned in Section II and shown in Figure 2, Graph A is the “remote” graph in the TableMult implementation. We see significant speedup when the extraction graph (E) is used in the first position, because fewer entries are sent over the network than when the E and G_{orig} switch positions.

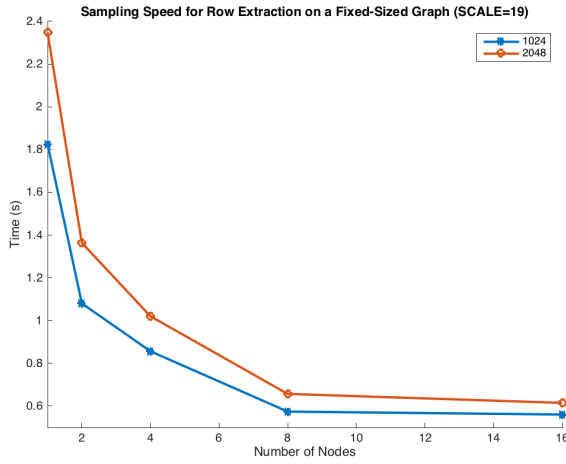


Fig. 4. Graphulo Row Extraction Experiment Results for Fixed Total Problem Size (Strong Scaling)

Figure 4 presents the results of our strong scaling experiment, where the graph size is $SCALE = 19$ and the number of extracted rows was 1024 or 2048 is kept fixed and we vary the number of processing elements. This graph clearly shows the expected drop-off in the time required to extract the number of rows as we increase the number of nodes. The sample times are within an acceptable range for interactive queries on a large dataset.

Figure 5 presents the results of our weak scaling experiment, where the graph $SCALE$ parameter changes from 15 to 19 (depending on the number of nodes) and the number of extracted rows is fixed to either 1024 or 2048. These results also indicate the expected constant-time performance as the amount of work per node is held constant. As before, we see a



Fig. 5. Graphulo Row Extraction Experiment Results for Fixed Problem Size Per Node (Weak Scaling)

performance difference of approximately 10% in computation time between the experiment on rows of 1024 and 2048.

These results demonstrate that Graphulo can provide a quick and effective method of extracting rows from a larger in-database graph for analysis. The extracted cued graph exists in the Accumulo database as a new table, available for additional local processing for analysis without anyone having to copy large amounts of data across the network or set up filtered iterators on the large graph. This capability enables the size of the cued graph data, not the size of the dataset from which it was sampled, to be the primary factor in Graphulo analytics.

V. DISCUSSION

Accumulo has many configuration parameters available for tweaking. The authors focused on scaling characteristics of Graphulo and provided a minimum amount of customization for our particular setup. A fully optimized Accumulo cluster should provide additional performance. However, we believe that the underlying scaling capabilities of Graphulo demonstrated above would only improve from additional tweaking.

Previous experiments used Accumulo’s ability to leverage the operating system’s native memory maps. Due to configuration management considerations, this option was not available on the shared cluster. Therefore, the above experiments were run using Java’s built-in memory management. Informal experimentation has demonstrated a 2x speedup in performance rate between the two options. Additional stability is possible by not relying on Java’s garbage collection system, which can sometimes result in timeouts and failed nodes across the entire cluster. Finally, the stability provided by native memory maps should provide opportunities for additional parallelization by enabling additional tablet splits per table.

At this time, column families, visibility labels, and timestamps are not taken into account in our Graphulo experiments. This work is potentially useful and interesting, especially for subgraph-based computations that may involve comparisons

TABLE II
NUMERICAL PARAMETERS AND RESULTS FOR THE GRAPHULO ROW EXTRACTION EXPERIMENTS

Nodes	SCALE	# Rows	Partial Prod	Seconds	Speedup	SCALE	# Rows	Partial Prods	Seconds	Speedup
1	19	1024	1.32×10^4	1.83	1.00	15	1024	1.51×10^4	0.69	1.00
		2048	2.67×10^4	2.35	1.00		2048	3.13×10^4	0.77	1.00
2	19	1024	1.41×10^4	1.08	1.80	16	1024	1.48×10^4	0.47	1.43
		2048	3.07×10^4	1.36	1.97		2048	2.98×10^4	0.53	1.38
4	19	1024	1.54×10^4	0.86	2.49	17	1024	1.54×10^4	0.56	1.24
		2048	3.07×10^4	1.02	2.63		2048	3.02×10^4	0.60	1.24
8	19	1024	1.54×10^4	0.57	3.71	18	1024	1.54×10^4	0.51	1.39
		2048	3.07×10^4	0.66	4.11		2048	3.09×10^4	0.53	1.44
16	19	1024	1.56×10^4	0.56	3.88	19	1024	1.56×10^4	0.56	1.26
		2048	3.02×10^4	0.61	4.29		2048	3.07×10^4	0.62	1.22

over data access policies, different time periods, or different types of graph information.

VI. CONCLUSIONS

In this paper, we presented experiments that demonstrated the scaling properties of Graphulo as its underlying Accumulo cluster grows from one node to sixteen nodes. Our first experiment was a large-scale graph multiplication task that extended previously published work. The second experiment demonstrated Graphulo's ability to quickly extract a set of vertices from a larger graph, which is useful in a more common use case where only a subset of a large stored graph is needed for analysis. In both experiments, we were able to demonstrate the expected strong and weak scaling in our processing rate. Additionally, we presented several areas of additional Graphulo improvement and experimentation.

For future work, we are interested in developing a greater suite of algorithms that can be directly developed using Graphulo such as dimensional analysis [20] and big data sampling [21]. Further, we would like to investigate providing a Pig interface to Graphulo as well as calling Graphulo operations from the Apache Spark framework.

ACKNOWLEDGMENTS

The authors would like to thank the Lincoln Laboratory Supercomputing Center team members: Bill Arcand, Bill Bergeron, David Bestor, Chansup Byun, Mike Houle, Matt Hubbell, Mike Jones, Anna Klein, Pete Michaleas, Lauren Milechin, Julie Mullen, Andy Prout, Tony Rosa, Siddharth Samsi, Chuck Yee, and Albert Reuther.

REFERENCES

- [1] J. Kepner, D. Bader, A. Buluç, J. Gilbert, T. Mattson, and H. Meyerhenke, "Graphs, matrices, and the GraphBLAS: Seven good reasons," *Procedia Computer Science*, vol. 51, pp. 2453–2462, 2015.
- [2] A. Buluç and J. R. Gilbert, "The combinatorial BLAS: Design, implementation, and applications," *International Journal of High Performance Computing Applications*, p. 1094342011403516, 2011.
- [3] V. Gadepally, J. Bolewski, D. Hook, D. Hutchison, B. Miller, and J. Kepner, "Graphulo: Linear algebra graph kernels for NoSQL databases," *CoRR*, vol. abs/1508.07372, 2015.
- [4] J. Kepner and J. Gilbert, *Graph Algorithms in the Language of Linear Algebra*. SIAM, 2011, vol. 22.
- [5] J. Kepner and V. Gadepally, "Adjacency matrices, incidence matrices, database schemas, and associative arrays," in *International Parallel & Distributed Processing Symposium Workshops (IPDPSW), 2014 IEEE International*. IEEE, 2014.
- [6] T. Mattson, D. Bader, J. Berry, A. Buluç, J. Dongarra, C. Faloutsos, J. Feo, J. Gilbert, J. Gonzalez, B. Hendrickson *et al.*, "Standards for graph algorithm primitives," *arXiv preprint arXiv:1408.0393*, 2014.
- [7] D. Bader, A. Buluç, J. G. UCB, J. Kepner, and T. Makson, "The graph BLAS effort and its implications for exascale," *SIAM Workshop on Exascale Applied Mathematics Challenges and Opportunities*, 2014.
- [8] R. S. Xin, J. E. Gonzalez, M. J. Franklin, and I. Stoica, "GraphX: A resilient distributed graph system on Spark," in *First International Workshop on Graph Data Management Experiences and Systems*. ACM, 2013, p. 2.
- [9] G. Malewicz, M. H. Austern, A. J. Bik, J. C. Dehnert, I. Horn, N. Leiser, and G. Czajkowski, "Pregel: a system for large-scale graph processing," in *Proceedings of the 2010 ACM SIGMOD International Conference on Management of data*. ACM, 2010, pp. 135–146.
- [10] "GovernmentCommunicationsHeadquarters/Gaffer." [Online]. Available: <https://github.com/GovernmentCommunicationsHeadquarters/Gaffer>
- [11] Y. Low, D. Bickson, J. Gonzalez, C. Guestrin, A. Kyrola, and J. M. Hellerstein, "Distributed GraphLab: A framework for machine learning and data mining in the cloud," *Proceedings of the VLDB Endowment*, vol. 5, no. 8, pp. 716–727, 2012.
- [12] R. Angles and C. Gutierrez, "Survey of graph database models," *ACM Computing Surveys (CSUR)*, vol. 40, no. 1, p. 1, 2008.
- [13] D. Hutchison, J. Kepner, V. Gadepally, and A. Fuchs, "Graphulo implementation of server-side sparse matrix multiply in the Accumulo database," in *IEEE High Performance Extreme Computing Conference (HPEC)*. IEEE, 2015.
- [14] R. Sen, A. Farris, and P. Guerra, "Benchmarking Apache Accumulo bigdata distributed table store using its continuous test suite," in *IEEE International Congress on Big Data*. IEEE, 2013.
- [15] J. Kepner, W. Arcand, D. Bestor, B. Bergeron, C. Byun, V. Gadepally, M. Hubbell, P. Michaleas, J. Mullen, A. Prout *et al.*, "Achieving 100,000,000 database inserts per second using Accumulo and D4M," *IEEE High Performance Extreme Computing Conference (HPEC)*, 2014.
- [16] J. Kepner, J. Chaidez, V. Gadepally, and H. Jansen, "Associative arrays: Unified mathematics for spreadsheets, databases, matrices, and graphs," in *Proceedings of New England Database Day (NEDB)*, 2015.
- [17] D. Bader, K. Madduri, J. Gilbert, V. Shah, J. Kepner, T. Meuse, and A. Krishnamurthy, "Designing scalable synthetic compact applications for benchmarking high productivity computing systems," *Cyberinfrastructure Technology Watch*, vol. 2, pp. 1–10, 2006.
- [18] V. Gadepally and J. Kepner, "Using a power law distribution to describe big data," in *High Performance Extreme Computing Conference (HPEC)*. IEEE, 2015, pp. 1–5.
- [19] P. Burkhardt and C. A. Waring, "A cloud-based approach to big graphs," in *IEEE High Performance Extreme Computing Conference (HPEC)*. IEEE, 2015.
- [20] V. Gadepally and J. Kepner, "Big data dimensional analysis," *IEEE High Performance Extreme Computing*, 2014.
- [21] V. Gadepally, T. Herr, L. Johnson, L. Milechin, M. Milosavljevic, and B. A. Miller, "Sampling operations on big data," in *Proceedings of the 49th Asilomar Conference on Signals, Systems and Computers*. IEEE, 2015, pp. 1515–1519.