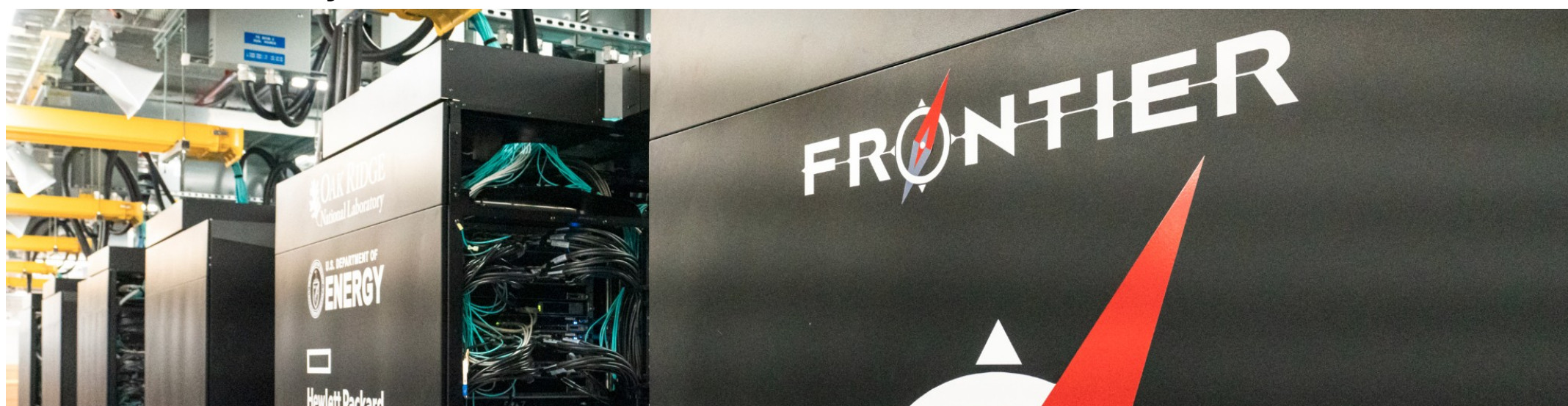


Algorithm-Based Fault Tolerance at Scale

Hayden Estes, Dr. Joshua Dennis Booth
Department of Computer Science

Overview

The rising need for fault tolerant systems is higher than ever due to the introduction of exascale computing. A system's fault tolerance is defined by its hardware, algorithms, and data types. Modern heterogeneous systems are composed of a large number of processors that support a variety of these data types. Understanding these data types' roles within an algorithm is vital to measuring fault tolerance. Exploring various floating-point formats shows how they can impact fault tolerance methods used for the sparse conjugate gradient algorithm (CG). We have determined that the modern Brain Floating Point Format (BFloat16) used in many Google tensor processing units (TPUs) can not be applied successfully, while IEEE754 32-bit floating point format (IEEE 32-bit) is utilized and optimized in GPUs (e.g. NVIDIA Tensor Core) can be applied successfully. These findings will make long-running scientific codes that use CG as a solver method able to ensure accuracy with minimal increased run time.



Pictured: Frontier, the world's first exascale supercomputer. Frontier will experience several faults every day due to its size and complexity,

Background

1. What are faults?

A fault is when a computer performs an operation incorrectly for various reasons.

- Hard-faults: Replicable physical faults in computer hardware.
- Soft-fault: transitive temporary faults. (e.g., alpha particle, packing pollution, cosmic radiation, etc)

2. How frequent are faults?

The frequency of faults is directly correlated to the complexity of the system.

- Faults on an average computer are reasonably rare due to low complexity.
- High-performance systems have several millions of processors which heavily increases the frequency of faults.

3. How can a system be more fault tolerant?

Faults are tolerated through use of detection, duplication, and correction.

- Detection involves checking the results in real time based on algorithms or generalized methods regardless of the algorithm. (e.g. duplicating calculations)
- Correction involves accepting the majority vote, continuing from a checkpoint, etc.

Key References

- Manu Shantharam, Sowmyalatha Srinivasamurthy, and Padma Raghavan. 2012. Fault tolerant preconditioned conjugate gradient for sparse linear system solution. In Proceedings of the 26th ACM international conference on Supercomputing (ICS '12). Association for Computing Machinery, New York, NY, USA, 69–78.
- Berrocal, Eduardo et al. "Lightweight Silent Data Corruption Detection Based on Runtime Data Analysis for HPC Applications." *Proceedings of the 24th International Symposium on High-Performance Parallel and Distributed Computing*. ACM, 2015. 275–278. Web.
- Cappello, Franck et al. "Toward Exascale Resilience." *The international journal of high performance computing applications* 23.4 (2009): 374–388. Web.
- Sun, Hongyang et al. "Selective Protection for Sparse Iterative Solvers to Reduce the Resilience Overhead." 2020 IEEE 32nd International Symposium on Computer Architecture and High Performance Computing (SBAC-PAD). IEEE, 2020. 141–148. Web.

Acknowledgments

Thank you to the *National Science Foundation* for funding this project through NSF Grant 2135310, the *Extreme Science and Engineering Discovery Environment* for access to high performance hardware. The RCEU program is sponsored in part by the UAH Office of the President, Office of the Provost, Office of the Vice President for Research and Economic Development, The Dean of the College of Science, the Dean of the College of Engineering, and the Alabama Space Grant Consortium.

Methodology

Continuing the work of [1], we decided to explore how limiting levels of precision may affect the efficiency of fault tolerance. We implemented our own preconditioned CG algorithm so that we had full control over the precision during every mathematical operation. Using this implementation, we supplied a set of SuiteSparse matrices and applied the data types IEEE 64-bit, IEEE 32-bit, and BFloat16. After confirming our data was correct through duplication and comparing the vectors of the previous work, we were able to analyze our results for any data types' benefits.

Results

Ultimately, we determined that the CG iterations, memory overhead, and overall precision are the key aspects of these data types. Figures 1 and 2 below show characterizing results of our preconditioned CG algorithm.

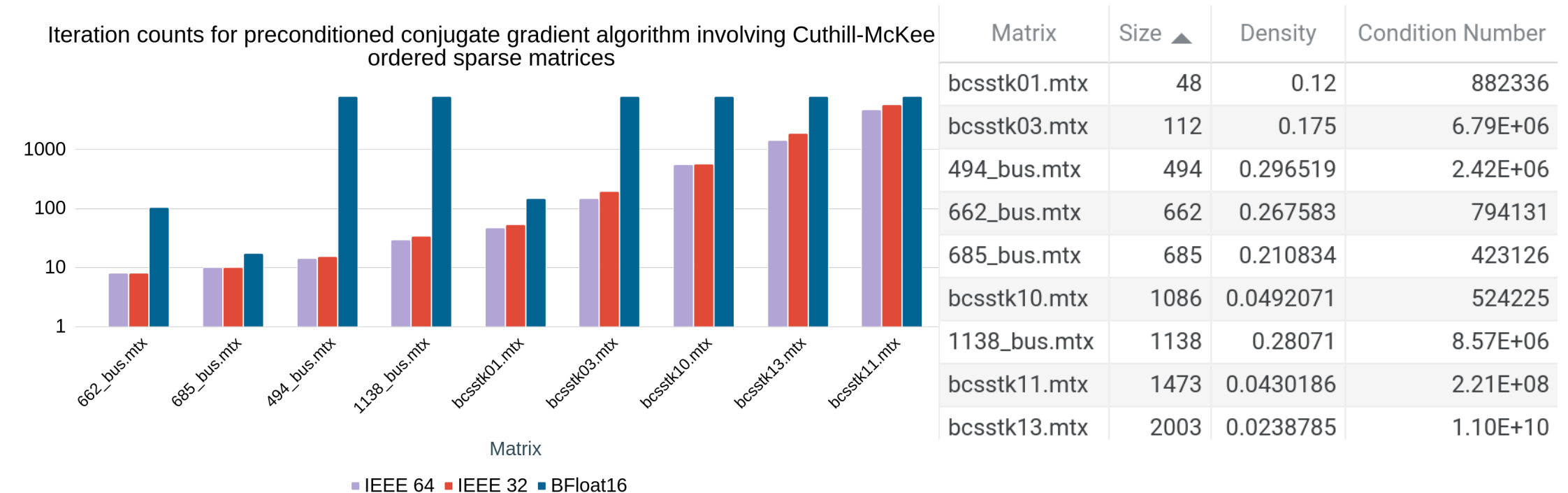


Figure 1: PCG iteration counts with varying precisions
*Note: The iterations have an artificial maximum of 8000

Figure 2: Sparse Matrices with their respective size, density, and condition values

IEEE 32	vs	BFloat16
→ More precise decimals		→ Less precise decimals
→ Large iteration decrease		→ Large iteration increase
→ More memory overhead		→ Less memory overhead
IEEE 32	vs	IEEE64
→ Less precise data type		→ More precise data type
→ Minimal iteration increase		→ Large iteration increase
→ Less memory overhead		→ Less memory overhead

- IEEE 32-bit precision is a beneficial alternative to IEEE 64-bit for CG algorithms while still causing a minimal increase in iteration count due to the efficiency increase in fault tolerance methods.
- BFloat16 is not a viable alternative to IEEE 64-bit or IEEE 32-bit for CG algorithms due to the large increase in iteration counts compared to the fault tolerance efficiency

Impact

Our results make CG more fault tolerance feasible due to the higher duplication efficiency with the IEEE 32-bit data type, in particular allowing GPUs to be used. CG algorithms are used to solve sparse systems of linear equations, so they are included in many other applications represented by partial differential equations. All long-running programs utilizing CG algorithms such as climate, electrical circuit, and robotic manipulator simulations can be improved using our findings.