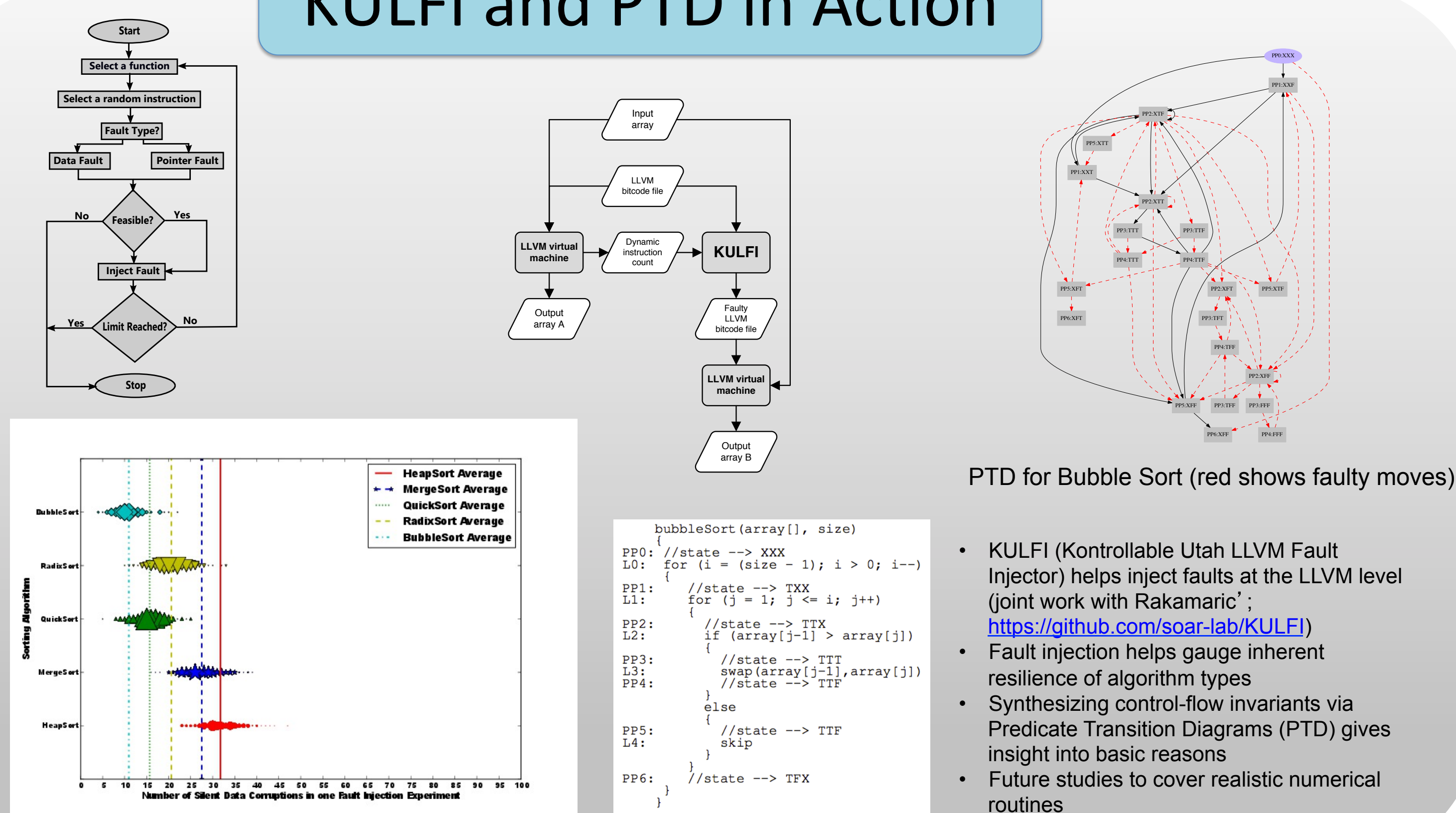
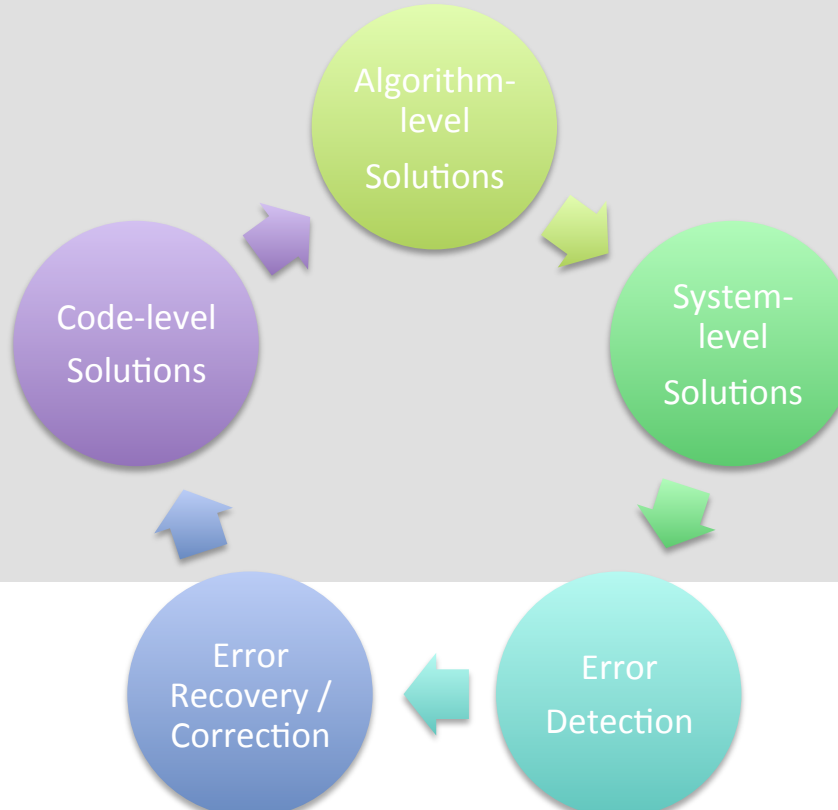


KULFI and PTD in Action



Abstract

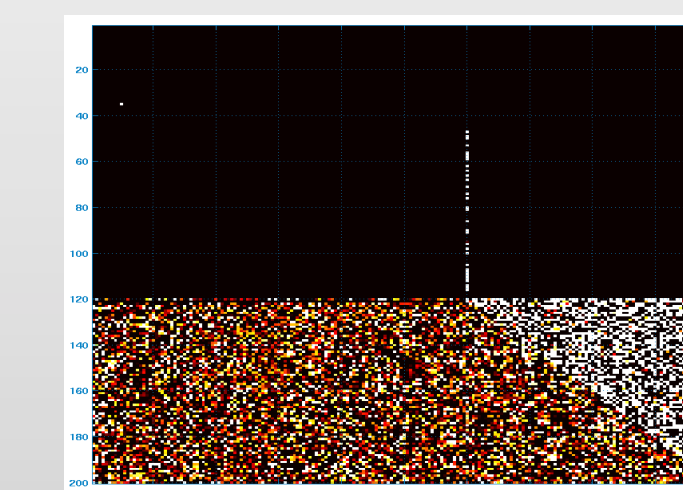
We are addressing the problems in software resilience with a holistic multifaceted approach that spans across software levels. One approach emphasizes tracking expected control flows or data invariants, and is aimed at detecting silent data corruption. Another explores language extensions and compiler technology to convey to compilers and run-time system resilience properties of code sections and algorithms. Additionally, we are investigating specific algorithmic properties of applications to develop fault tolerant extensions to dense and sparse methods. At the highest level, we detect silent-data corruptions by replicating and comparing values across MPI processes.



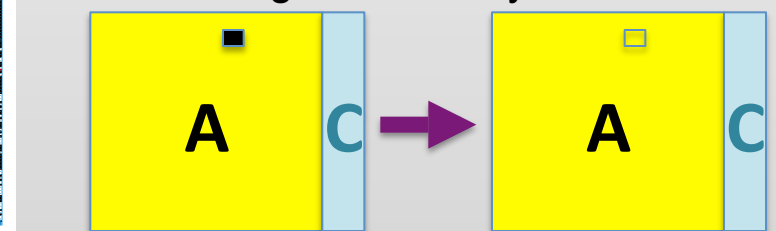
Algorithmic-Based Approaches

CHALLENGE

Detection and Correction of Errors using algorithmic properties



Principle of Algorithm Based Fault Tolerance: the data is augmented with checksum columns that preserve a mathematical property, introducing redundancy in the data.



The figure illustrates how a single memory corruption propagates to the entire matrix during a LU factorization due to the recursive nature of the algorithm.

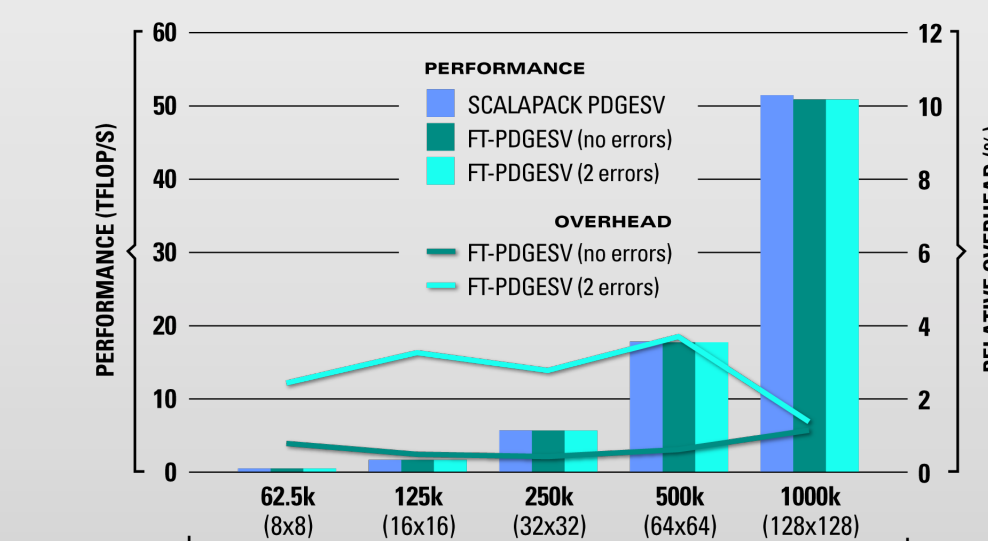
In case of error (e.g. memory corruption), checksum inversion allows to restore the corrupted value.

Scope of ABFT: although ABFT techniques have been applied mostly for dense direct method in linear algebra (one-sided factorizations, two-sided factorizations, ...), it is not limited to these: iterative methods (CG, ...) and sparse operations are also target of ABFT techniques.

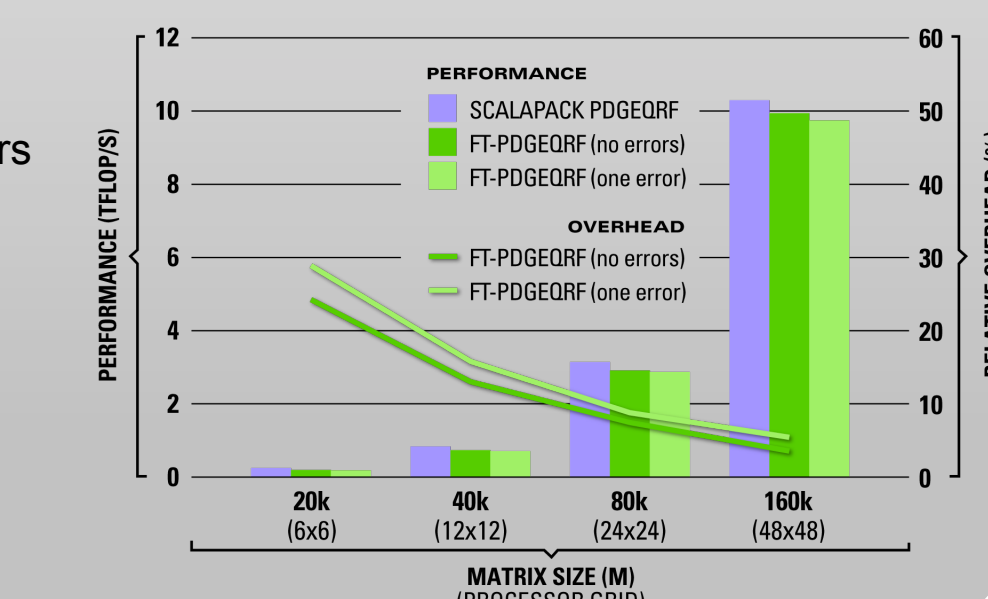
PERFORMANCE

Kraken – Weak Scaling

LU Factorization (PDGESV) Despite Memory Corruption Errors

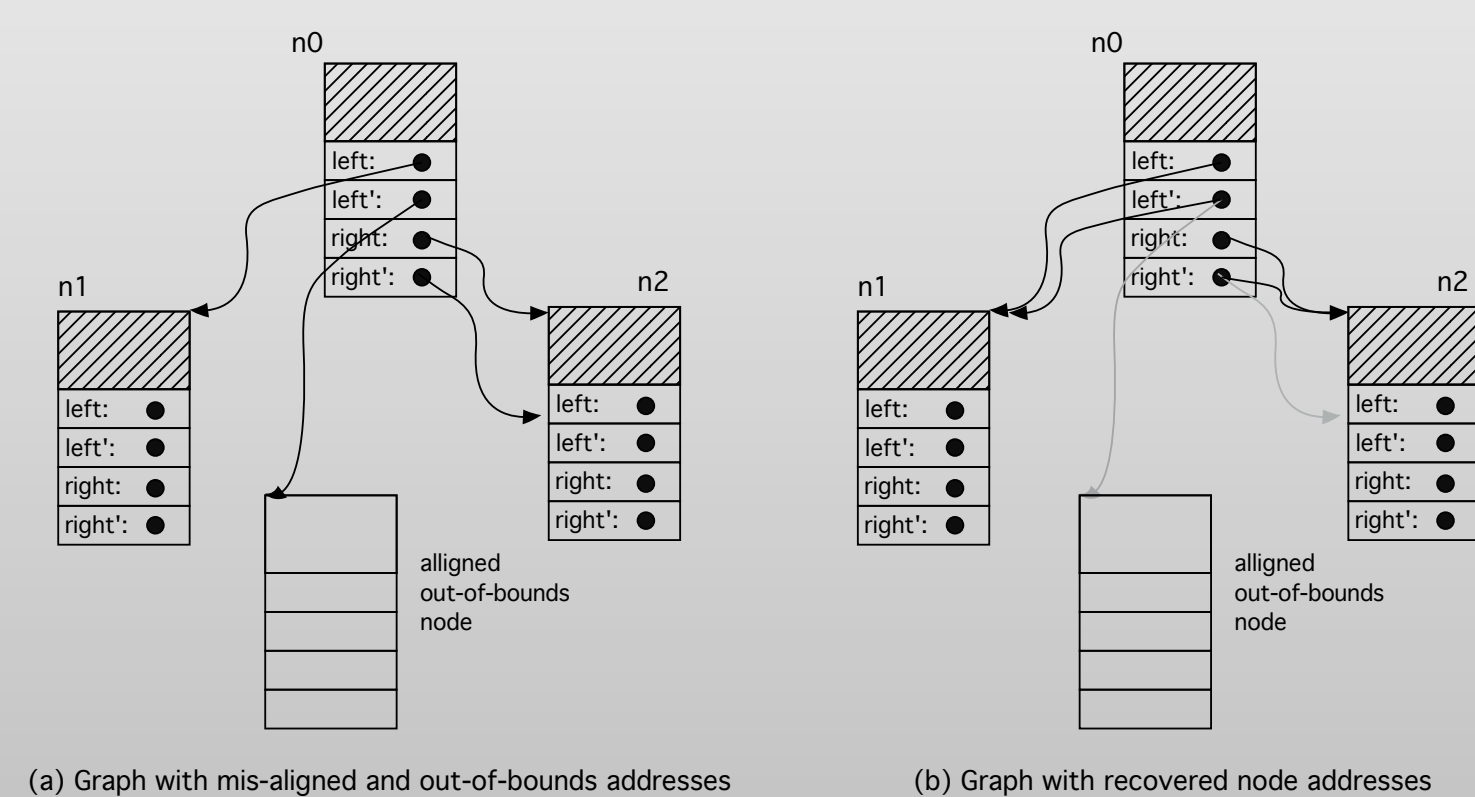
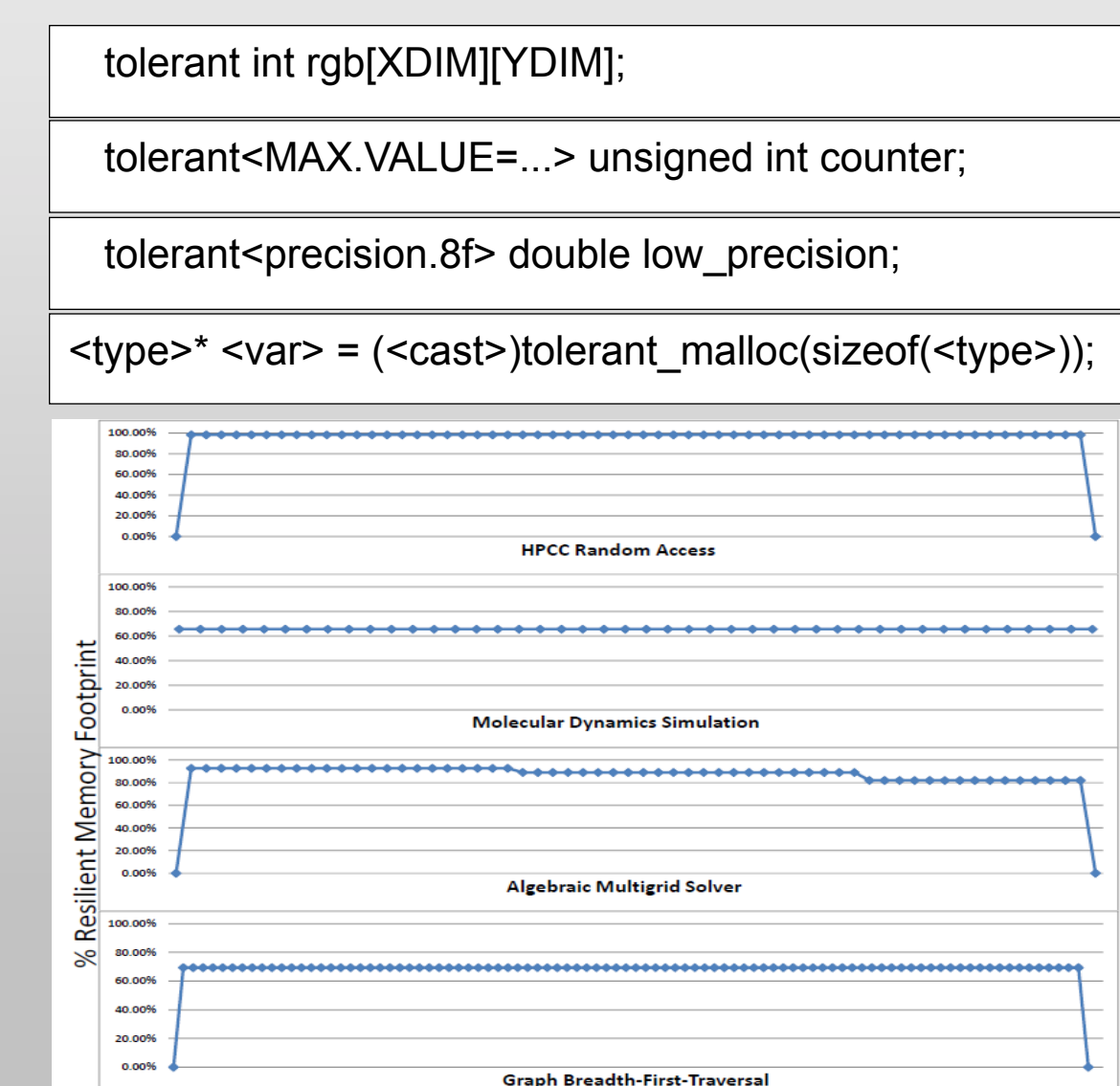


QR Factorization (PDGEQRF) Despite Fail-Stop Errors



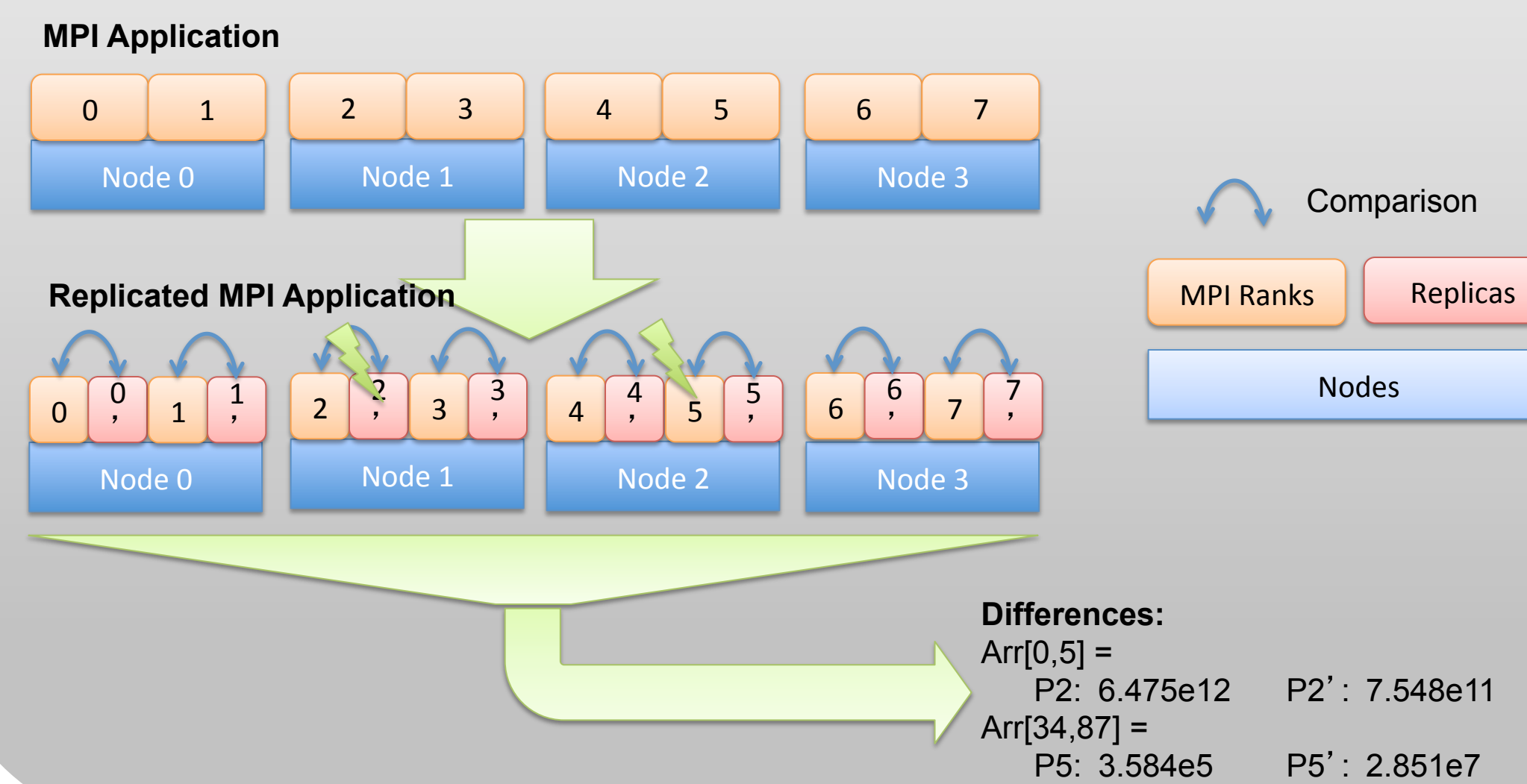
Language Extensions & Compiler Technology for Resilience

- Annotations that allow user to express fault-tolerant requirements and expectations: when and where errors matter and what to do about them.
- ROSE source-level resilience-oriented and user-guided transformations for array-based pointers and graph-based computations



Silent Error Detection

- We have seen unexpected behavior in jobs at scale on the LLNL Sequoia machine
 - Certain high-performance LINPACK runs have high residual
- Currently have no way to detect silent memory corruption
 - Conducting a detailed characterization of memory error rate of BG/Q, Cray
- Developing a tool, Dagnet, that finds memory errors through MPI replication
 - Replicate MPI processes on-node, do shared-memory comparison of arrays
 - Can convert any MPI program into a silent error detector



Comprehensive Algorithmic Resilience

We demonstrate that is possible to protect scientific applications from many sources of errors combining three different resilience techniques:

- Algorithmic error checkers
- Replications of key data structures
- Checkpoint-restart mechanisms

Significant improvements can be achieved in terms of reduction of performance slowdown and output accuracy

