Resilience Extensions for MPI: ULFM

User Level Failure Mitigation is a set of MPI interface extensions enabling Message Passing programs to restore MPI communication capabilities affected by process failures. It supports rebuilding communicators, RMA windows and I/O Files.

FLEXIBILITY

- No particular recovery model is imposed or favored. Instead, a set of versatile APIs is included that provides support for different recovery styles (checkpoint, ABFT, iterative, Master-Worker, etc.).
- Application directs the recovery, it pays only for the level of protection it needs.
- Recovery can be restricted to a subgroup, preserving scalability and easing the composition of libraries.

PERFORMANCE

- Protective actions are outside of critical MPI routines.
- MPI implementors can uphold communication, collective, one-sided and I/O management algorithms unmodified.
- Encourages programs to be reactive to failures, cost manifests only at recovery.

PRODUCTIVITY

- Backward compatible with legacy, fragile applications.
- Simple and familiar concepts to repair MPI.
- Portability guaranteed by standardization.
- Provides key MPI concepts to enable FT support from library, runtime and language extensions.

The ULFM specification is a crucial infrastructure that will enable the deployment of advanced, production quality fault tolerant techniques; it is a versatile solution to improve the efficiency of novel and established techniques.



Resilience Extensions for MPI: ULFM

ULFM provides targeted interfaces to empower recovery strategies with adequate options to restore communication capabilities and global consistency, at the necessary levels only.

CONTINUE ACROSS ERRORS

In ULFM, failures do not alter the state of MPI communicators. Point-to-point operations can continue undisturbed between non-faulty processes. ULFM imposes no recovery cost on simple communication patterns that can proceed despite failures.

Consistent reporting of failures would add an unacceptable

an operation is disrupted; other ranks may still complete their

operations. A process can use MPI [Comm,Win,File] revoke to

propagate an error notification on the entire group, and could, for example, interrupt other ranks to join a coordinated recovery.

performance penalty. In ULFM, errors are raised only at ranks where



COLLECTIVE OPERATIONS

GROUP EXCEPTIONS

Allowing collective operations to operate on damaged MPI objects (Communicators, RMA windows or Files) would incur unacceptable overhead. The MPI_Comm_shrink routine builds a replacement communicator, excluding failed processes, which can be used to resume collective communications, spawn replacement processes, and rebuild RMA Windows and Files.



OPEN MPI ULFM IMPLEMENTATION PERFORMANCE



Sequoia AMG is an unstructured physics mesh application with a complex communication pattern that employs both point-to-point and collective operations. Its failure free performance is unchanged whether it is deployed with ULFM or normal Open MPI.





The failure of rank 3 is detected and managed by rank 2 during the 512 bytes message test. The connectivity and bandwidth between rank 0 and rank 1 are unaffected by failure handling activities at rank 2.