

PMIx Fault Events Detection and Dissemination use case: Fault-Tolerance in Open MPI 5 and more

Aurelien Bouteiller, Zhong Dong, George Bosilca



ICL
INNOVATIVE
COMPUTING LABORATORY

T
THE UNIVERSITY OF
TENNESSEE
KNOXVILLE

PMIx Fault Events Interface

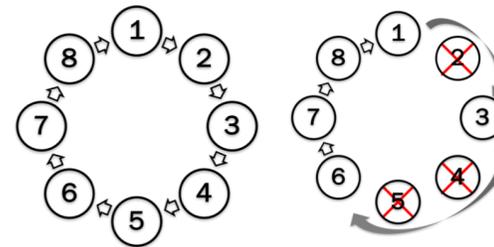
Dong Zhong, Aurelien Bouteiller, Xi Luo, and George Bosilca. 2019. **Runtime Level Failure Detection and Propagation in HPC Systems**. In 26th European MPI Users' Group Meeting (EuroMPI 2019), September 11–13, 2019, Zürich, Switzerland. <https://doi.org/10.1145/3343211.3343225>

FRONT END: PMIx Interface

- Control of the detection service through PMIx interfaces
- Fault Events are presented as normal PMIx Events**
- Very simple to hook into
- Compatible with non-ft builds (event is just not generated)
 - `PMIX_Register_event_handler(..., PMIX_ERR_PROC_ABORTED, ...)`
- Opens the gate for efficient management of failures for an emerging field of libraries, programming models, and runtime systems operating on large-scale systems.
- Well specified interface helps **usability by multiple client types**

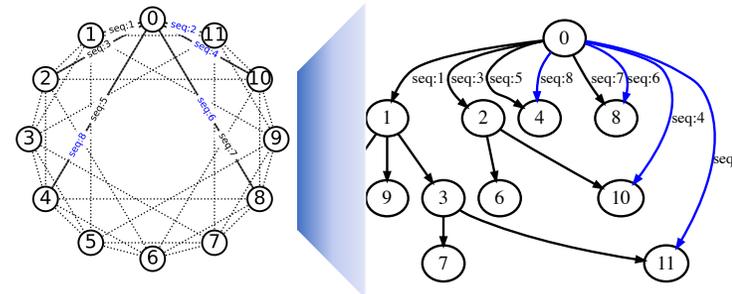
BACK END Capability integrated in **PRTE 2.0**, other LM/SMS can produce the same events using their own methodologies

FAILURE DETECTION COMPONENT



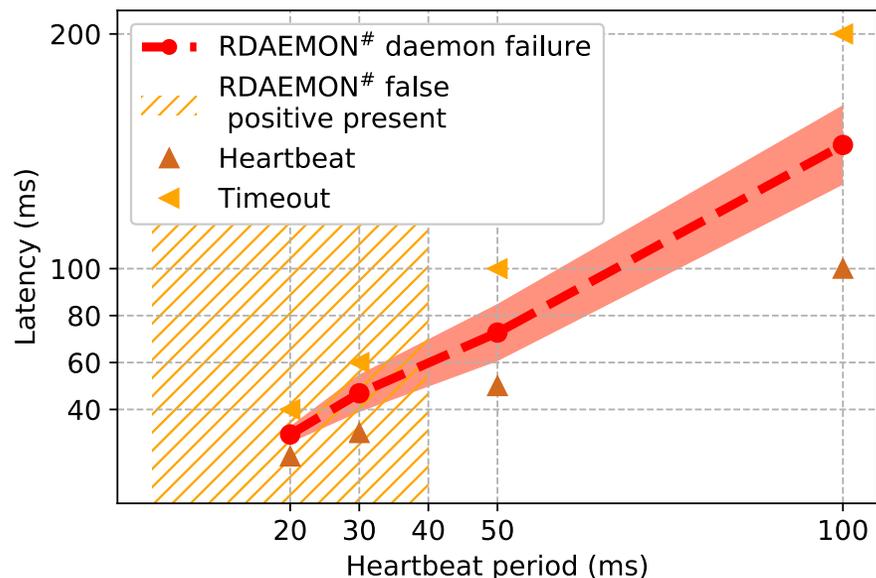
Node/PRTE servers failure detected with **ping on a mendable ring**
Minimal ping noise

RELIABLE EVENT PROPAGATION COMPONENT



Event Propagation along **multiple binomial trees** extracted from a circulant graph: **fixed degree (log), fast and resilient**

Advantages in delegating detection to the PMIx runtime



- Accuracy of detection is very good (in the order of 100ms can be achieved in practice at scale)
- False detection rate **independent of the application communication pattern**
 - *Prior MPI-based detector would produce false positive when application does not call MPI procedures*
- *Reusable in different programming models*

Performance variability in GRAPH500 with an active PMIx-PRRTE Failure detector
 Left: MPI Right: OpenSHMEM

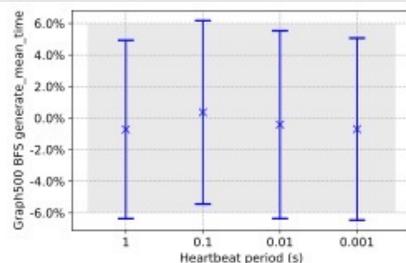


Figure 17: Overhead for generating BFS running `mpi.test.simple` when using PRRTE with fault tolerance over PRRTE (32K MPI ranks; the gray area represents the normal variability of the benchmark).

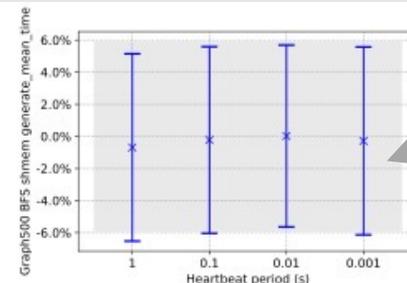


Figure 19: Overhead for generating BFS running `graph500.shmem.one.sided` upon PRRTE with fault tolerance over PRRTE (32K OpenSHMEM PEs; the gray area represents the normal variability of the benchmark).

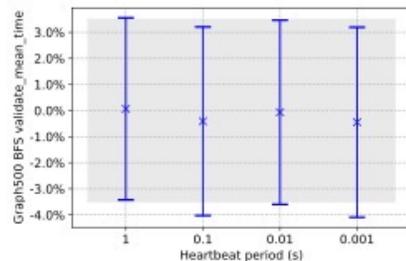


Figure 18: Overhead for validating BFS in `mpi.test.simple` when using PRRTE with fault tolerance over PRRTE (32K MPI ranks; the gray area represents the normal variability of the benchmark).

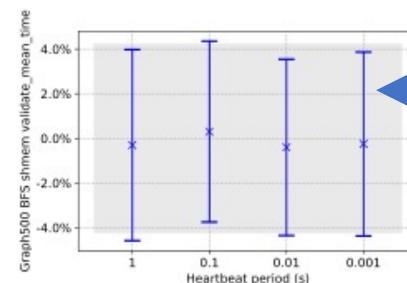


Figure 20: Overhead for validating BFS running `graph500.shmem.one.sided` upon PRRTE with fault tolerance over PRRTE (32K OpenSHMEM PEs; the gray area represents the normal variability of the benchmark).

Gray area represents normal benchmark variability

Blue error bars show the variability as measured with detection ON

Experiments performed on NERSC's Cori: Cray XC40 supercomputer with Intel Xeon "Haswell" processors and the Cray "Aries" high speed inter-node network, 32 cores per node, 32K processes total.



THE UNIVERSITY OF
 TENNESSEE
 KNOXVILLE

Fault Management in PMIx

What works, and Future Directions

- PMIx Event Handling proved to be a natural fit to disseminate fault events
- PMIx abstracts the ‘detector’, which lets application write portable fault-tolerant code that can operate on a variety of PMIx servers (including non-FT ones)
- PMIx servers have freedom of implementation behind the API curtain

Future Extensions

- Fault-events implemented only in PRTE at the moment (PRTE can be launched under Slurm/Jsrn/ALPS, etc., in managed mode, but not all have native support yet).
- Current enabling of fault-detection/management is command-line based
- Can we move to a programmatic way to turn-on/turn-off resilient features?
- PRTE: Resilient overlay communication for commands/modex (reduce vulnerability to startup faults)
- Could user have fine-grain control over what operation (e.g., PMIx modex or Fence) are resilient, or not (performance optimization)?

The left half of the image features a background of a circuit board with various components and traces. Overlaid on this is the ICL logo, which consists of the letters 'ICL' in a bold, black, sans-serif font. The letter 'I' is solid black, while the 'C' and 'L' have an orange accent on their right and bottom edges, respectively. Below the logo, the words 'INNOVATIVE' and 'COMPUTING LABORATORY' are stacked in a smaller, black, sans-serif font.

ICL
INNOVATIVE
COMPUTING LABORATORY

The right half of the image shows a grayscale photograph of a large, multi-story brick building with a prominent clock tower. The building has Gothic-style architectural elements like pointed arches. In the foreground, there are trees and a street lamp. A faint grid pattern is overlaid on the entire image.

THE UNIVERSITY OF
TENNESSEE
KNOXVILLE