

Application-driven Fault-Tolerance for High Performance Distributed Computing.

Franck Cappello Argonne/MCS, UIUC

George Bosilca University of Tennessee at Knoxville





Agenda

• Introduction: 0h30

- 5 minutes motivation (failure rates) Franck
- 10 minutes checkpoint/restart (application level versus system level) Franck
- 15 minutes alternatives (ULFM, Replication, etc.) George.

• VeloC and Hands-on installation: 1h15 - Franck

- 25 mins -> VeloC presentation
- 5 mins -> questions
- 45 mins -> hands-on-session
 - 15 minutes of buddle installation
 - 30 minutes of VeloC (Bogdan on rescue, if needed)
 - Configuration: 10 minutes
 - Filling the gap in the heat equation: 10 minutes
 - Playing with failures: 10 minutes

• ULFM: 1h15 – George

- 25 mins -> ULFM presentation
- 5 mins -> questions
- 45 mins -> hands-on-session
 - 30 min on ULFM
 - 15 minutes on ULFM + VeloC.

Why FT for HPC?

- FT need for HPC was marginal because HPC system MTBF were high enough (1 week, 1 month). This is not true anymore for large systems today (MTBF of 1 day and less are seen)
- It can only get worse with the increase of the number of components and component complexity
- There is no compromise:
 - Fault tolerance is not like other problems of HPC (performance, efficiency, power consumption, etc.) \rightarrow there is no half success:
 - Application execution succeeds with correct results or fails!
- Clouds are starting considering HPC applications and Cloud nodes have typically a much lower MTBF than HPC nodes.

Some important definitions

From Avizienis, Laprie et al.:

Definition from the notion of service: a sequence of the system's external states (perceived by users)

Correct service: is delivered when the service implements the system Example:

- Se A) A particle hits a DRAM cell and generates a fault
 - B) The fault changes the DRAM cell state and becomes an error
 - C) The error does not affect the rest of the system until a process reads
 - the cell
 - D) The error propagates as a failure if after the read of the memory cell the software computation, control or I/O deviates from the behavior it would have had from a correct memory cell
- Fault: The adjudged or hypothesized cause of an error (*root cause of the failure*)

Specific Error Outcomes in HPC

Types of errors:

- Power outage
- Hard errors (broken component: memory, network, core, disk, etc.)

Detected coff errors (hit flip in memory

To tolerate these errors, current systems essentially follow a masking approach:

-Mask transient/soft errors at the hardware level

-Checkpoint-restart for process failure (errors that eventually lead to application crash)

• User errors (Human)

cause by system type.

Hardware Software

Jetwork

Environmen Human

G. Gibson

100

80

Classes of errors:

- Detected and corrected (by ECC, Replication, Re-execution)
- Detected and uncorrectable (leading to application crash)
- Undetected (leading to data corruption, application hang, etc.)

MTBF 10 Years Ago?

IBM Research

Blue Gene Hardware Reliability: Argonne Data

- BG/L System Design Target:
 - 64 Racks/131k cores MTBF should be greater than 7 Days

Paul Coteus, IBM

- Comparison of actual data made by ANL Labs
 - asked a number of facilities for reliability data.
- Multi teraflop IA64 or X86 systems have 100's to 1000's of individual compute nodes.
- For comparison between different systems, fail rates are normalized to peak system performance in teraflops

System	Peak System	Full System	Failures	Failures per	Failures
Туре	Performance	Mean Time	per Month	Month per	per Month
	(Teraflops)	Between Failures		Teraflop	per BG
		(Days)			Rack
IA64	3.0	1.3	24.0	8.000	
IA64	10.7	1.1	28.3	2.645	
x86	1.7	4.5	6.7	3.941	
x86	17.2	0.7	45.1	2.622	
Power 5	15.0	1.1	19.0	1.267	
Blue Gene	365.0	7.5	4.0	0.011	0.06

MTBF 5 years ago

Two classes:

Based on proprietary components: IBM designs BG line with a full system
 MTBF of 7 days (true for BG/L, BG/P, BG/Q?)
 Paul Coteus, IBM

	FIT per	Components per	FIT per
Component	Component	64K System	System
DRAM	5	608,256	3,041 K
Compute + I/O ASIC	20	66,560	1,331K
ETH Complex	160	3,024	484K
Non-redundant power supply	500	384	384K
Link ASIC	25	3,072	77K
Clock chip	6.5	1,200	8K
Total FITs			5,315K

Table 6.12: BlueGene FIT budget.

- Using commodity components (Intel, AMD processors, etc.): MTBF of about 1 day (some less, some more) for systems with 100,000+ cores
 - CEA Tera100 (6 in the top500) 20h MTBF for the whole machine, 4300 nodes, 140 000 cores, 500 Gb/s global file system bandwidth, 300 TB memory
 - Jeffrey Vetter (ORNL):

Stable

Jaguar XT5 status, April 2009

1: 32 hours

- Driver for downtimes: Spider testing
- · System has been up as long as 10 days

Current failure rates

Fault	Local Consequence	Cascading Consequence	Mean Time
			between Faults
Node failures	User processes running on	Full user execution crash	BW <mark>3</mark> 6.7h
(some hardware	the node crashes	because the runtime of the	[Mar14]
or OS part of the		resource/job manager decides	Titan ³ : 7.5 h
node fail leading		to kill the execution (R1) or	[Tiw14]
to a complete		because of a cascading to full	resulting in
failure of the		system outage (R2)	mean time to
node) ²			application
			failure of 40 h
			[TGR15]
Network failure	The user processes that	Potential full user execution	BW: 20 h (link
	cannot communicate	crash because of R1. Also if	failure)
	experience time-outs on	the execution was not able to	[Mar14]
	communication. OS or	checkpoint because of	
	runtime may kill these	network failure, then it will	
	processes. The affected	need to restart from the	
	processes may crash on	previous checkpoint (C1).	
	their own. However, user		
	processes may be able to		
	tolerate transient network		
	shoot down/rerouting.		

2 For example, GPU bus errors (disconnection of the GPU), voltage fault, kernel panic, PCI width degrade, machine check exception, and SXM (PCI) power off observed in Titan lead to process crashes [Gup15].

3 Time between failures of any node in the system. Each node MTBF is typically 25 years in these systems [Tiw14].

Interval between failure can be << MTBF



Devesh Tiwari, Saurabh Gupta, Sudharshan Vazhkudai, Lazy Checkpointing: Exploiting Temporal Locality in Failures to Mitigate Checkpointing Overheads on Extreme-Scale Systems, Proceedings of the Annual IEEE/IFIP Int'l Conference on Dependable Systems and Networks (DSN), 2014.

Documents and tools

- Fault tolerance for Distributed system is a not a new but it is a young domain for Parallel Commuting
- Books/articles related to Faults, Errors, Failures and distributed computing
 - A. Avizienis et al. "Basic Concepts and Taxonomy of Dependable and Secure Computing", IEEE Transactions on dependable and secure computing, Vol.1, No 1 January-March 2004
 - N. Lynch "Distributed Algorithms", Morgan Kaufmann Publishers Inc. 1996 ISBN:1558603484
 - E. Elnozahi "A Survey of Rollback-Recovery Protocols in Message Passing Systems", ACM Computing Survey, Vol. 34, No. 3, pp. 375-408, September 2002."
- Article related to FT/Resilience in Parallel Computing
 - F. Cappello et Al. "Toward Exascale Resilience". IJHPCA 23(4): 374-388 (2009)
 - F. Cappello "Fault Tolerance in Petascale/ Exascale Systems: Current Knowledge, Challenges and Research Opportunities". IJHPCA 23(3): 212-226 (2009)
 - M. Snir, R. W Wisniewski, J. A Abraham, S. V Adve, S. Bagchi, P. Balaji, J. Belak, Pradip Bose, F. Cappello, B. Carlson, A. A Chien, P. Coteus, N. A DeBardeleben, P. C Diniz, C. Engelmann, M. Erez, S. Fazzari, A. Geist, R. Gupta, F. Johnson, S. Krishnamoorthy, S. Leyffer, D. Liberty, S. Mitra, T. Munson, R. Schreiber, J. Stearley, E. Van Hensbergen, Addressing Failures in Exascale Computing, International Journal of High Performance Computing Applications, vol. 28, num. 2, pages 127-71, May 2014.
 - C. Di Martino, Z. Kalbarczyk, R. K. Iyer, F. Baccanico, J. Fullop, and W. Kramer. 2014. Lessons Learned from the Analysis of System Failures at Petascale: The Case of Blue Waters. In Proceedings of the 2014 44th Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN '14).
 - C. Di Martino, W. Kramer, Z. Kalbarczyk, R. Iyer, Measuring and Understanding Extreme-Scale Application Resilience: A Field Study of 5,000,000 HPC Application Runs, in Dependable Systems and Networks (DSN), 2015
 - D. Tiwari, S. Gupta and S. Vazhkudai, Lazy Checkpointing: Exploiting Temporal Locality in Failures to Mitigate Checkpointing Overheads on Extreme-Scale Systems, IEEE/IFIP International Conference on Dependable Systems and Networks, DSN 2014
 - D. Tiwari, S. Gupta, G. Gallarno, J. Rogers, and D. Maxwell. Reliability lessons learned from GPU experience with the Titan supercomputer at Oak Ridge leadership computing facility. In Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis (SC '15), 2015.
 - D. Tiwari, S. Gupta, J. Rogers, D. Maxwell, P. Rech, S. Vazhkudai, D. Oliveira, D. Londo, N. DeBardeleben, P. Navaux, L. Carro, A. Bland, A., Understanding GPU errors on large-scale HPC systems and the implications for system design and operation, High Performance Computer Architecture (HPCA), 2015
 - B. Schroeder and G. Gibson, Understanding failures in petascale computers. Journal of Physics: Conference Series 78:012022, 2007

Principle of Checkpoint restart

Full execution Checkpoint-Restart:



Checkpointing techniques

A checkpoint is just a 'snapshot' of a process (or system) at a certain point in time

A checkpointing system provides a way to take these snapshots, and to restart from them

Type of checkpoint mechanisms:

Kernel & User (System) Level Stack Easy to add checkpointing to existing code SP Works with (almost) any programs sbrk(0) General, 'coarse', approach (the full state of each process is saved) Heap &edata Examples: Libckpt, BLCR, DMTCP (Boston University) Data (Static) &etext Text

Application Level

Could require modifications Different API (memory level, file level) 'fine' grain approach: only the sections of the state necessary to restart is saved Examples: FTI, SCR, VeloC



Parallel Checkpointing

Coordinated Checkpoint

The objective is to checkpoint the application when there is no in transit messages between any two nodes Coordination:

 \rightarrow Automatic through a protocol (Chandy-Lamport) \rightarrow not needed in SPMD apps.

 \rightarrow Implicit: application level

Uncoordinated Checkpoint

No global synchronization (scalable) →Nodes may checkpoint at any time (independently of the others) →Need to log undeterministic events: Intransit Messages → too complex





Application Level Implicitly Coordinated Parallel Checkpointing

Example: heat distribution SPMD code

```
while(i < ITER_TIMES) {</pre>
        localerror = doWork(nbProcs, rank, M, nbLines, g, h);
        if (((i % ITER_OUT) == 0) && (rank == 0))
            printf("Step : %d, error = %f\n", i, globalerror);
        if ((i % REDUCE) == 0)
            MPI_Allreduce(&localerror, &globalerror, 1, MPI_DOUBLE, MPI_MAX,
MPI_COMM_WORLD);
        if (globalerror < PRECISION)
            break;
        i++;
        if (i % CKPT_FREQ == 0) {
             FILE *outFile = fopen("checkpoint", "wb");
             fwrite(&I, sizeof(int), 1, outFile);
             fwrite(h, sizeof(double), M * nbLines, outFile);
             fwrite(g, sizeof(double), M * nbLines, outFile);
        }
    }
```

Where to checkpoint

Example: SCR (LLNL) result: Aggregate checkpoint bandwidth to node-local storage scales linearly on Coastal



When to checkpoint

[Young 74] Let's assume that our system failure rates follow the bath tub pattern. We are interested to compute the checkpoint interval for the constant failure rate regime



Well modeled by the Exponential distribution Failure density function:

$$f(x;\lambda) = \begin{cases} \lambda e^{-\lambda x}, & x \ge 0, \\ 0, & x < 0. \end{cases}$$

MTBF = 1 / λ Failure rate

The main formula used to compute checkpoint Intervals in HPC systems.

A more accurate formula by John Daly that integrate restart time.

Interval = $\sqrt{2 \times checkpoint-time \times MTBF}$

John W. Young, « A first order approximation to the optimum checkpoint interval », Communications of the ACM, Volume 17 Issue 9, Sept. 1974

$$\tau = \sqrt{2\delta(M+R) \over _{\rm MTBF}} - \delta_{_{\rm Rst \ time}} - \delta_{_{\rm Ckpt \ time}}$$