

Motivation

Addressing exascale challenges:

- Resiliency, power, memory access, concurrency and heterogeneous hardware.
- Codes/algorithms need to be flexible and extensible to unknown architectures.
- Errors and faults will become more frequent compared to current platforms.

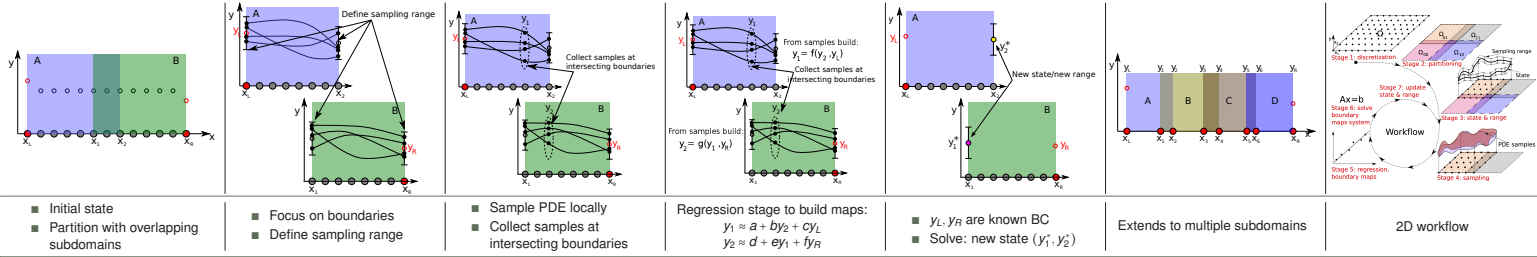
Objectives

- Accurate and robust simulation even in the presence of system faults.
- Assess predictive fidelity of extreme-scale scientific simulations.
- High performance on "uncertain" architectures.
- Portable, extensible and reusable codes.

Technical Approach

- What?
 - Resilient preconditioner for PDEs.
 - To achieve resiliency to both soft and hard faults.
- How?
 - Recast the original PDE as a sampling problem.
 - Focus solely on the information available.
 - State update through resilient data manipulation.

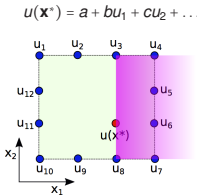
Resilient EXtreme Scale Scientific Simulation (REXSSS) Algorithm Overview



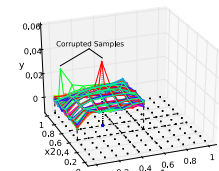
Regression stage: Laplace (ℓ_1) vs Gaussian (ℓ_2)

Regression Resilience

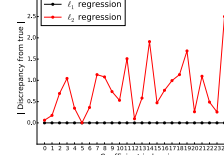
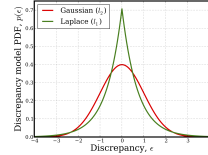
High-dimensional boundary maps:



Generate samples (some corrupted)
Run regression using collected samples

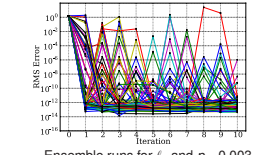
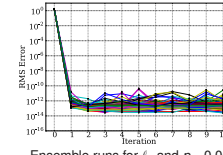


ℓ_1 : robust against presence of corrupted data
 ℓ_2 : highly sensitive to corrupted data



1D Elliptic PDE Test

- Bit-flips inserted during sampling with probability p .
- Laplace likelihood yields resiliency to soft faults.

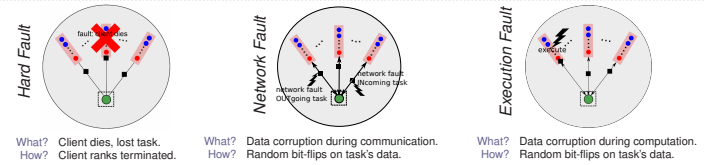


REXSSS Implementation

Server-Client Model (SCM)

- Cluster: 1 server + n clients.
- Servers:
 - Communicate between each other.
 - Safe data/state storage (sandboxed).
- Clients:
 - Independent from one another.
 - Only serve as computing units.
- Separates state from computation: reduces the overall vulnerability.
- Fault-tolerance via *User Level Fault Mitigation* ULFM-MPI (fault-tolerance.org).
- It aligns with the vision of future exascale architectures involving heterogeneous and hierarchical hardware required to meet energy and cost constraints.

Currently Supported Fault Models



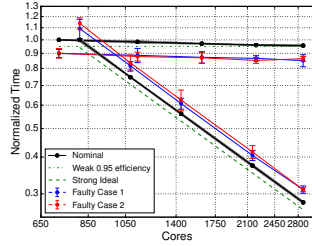
ULFM-MPI and SCM: Why is this a good combination?

- Server simply continues the execution using only the clients that are alive.
- Avoid ULFM-MPI collective procedures to rebuild the client/communicators.
- Servers probe the corresponding cluster communicator using `MP_I_ANY_SOURCE` to assess whether a new message is arriving from one of the clients.

Scalability

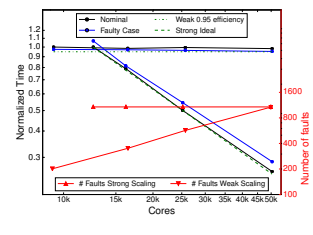
Uncertainty due to SDCs and client failures

- 1% failed clients, SDC=0.1% total tasks (—)
- 2% failed clients, SDC=0.1% total tasks (—)
- Times normalized by smallest nominal case: highlights scaling and overhead.
- Faulty runs show larger variability due to random loss of resources, and additional regression overhead to overcome silent data corruptions (SDCs).
- Red curve has larger overhead because the number of failed clients doubles.



Scaling on Edison (NERSC)

- Of total faults: SDCs=97%, ranks failed=3.0%.
- Constant machine fault rate.
- Times normalized by smallest nominal case: highlights scaling and overhead.
- Excellent scalability with and without faults:
 - within 95% for weak scaling.
 - within 90% for strong scaling.
- Overhead with respect to nominal case:
 - a downward shift for weak scaling.
 - an upward shift for strong scaling.



Energy and Resilience

Exploit resilience for energy purposes?

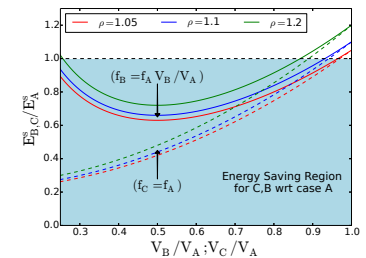
- Idea: lower the energy consumption during the *sampling stage* by means of voltage scaling.
- Compare three scenarios:
 - machine running at full operational capacity/speed
 - voltage/freq scaling on clients during sampling
 - voltage scaling on clients during sampling
- Same problem, SC configuration, and machine.
- Servers run at full capacity to keep the state safe.
- Power consumption (P) and energy (E) over time T :
 $P = \hat{P} + CV^2f$ $E = (\hat{P} + CV^2f)T$
frequency independent active power (\hat{P}), switch capacitance (C), voltage (V), frequency (f).

Energy Budget for Sampling Stage

- Full operational mode:**
 - $V_A, f_A, t_A^s =$ time for one task.
 - For N_A samples: $E_A^s = N_A(\hat{P}_A t_A^s + CV_A^2 f_A t_A^s)$
 - Reduced voltage/frequency:**
 - $V_B < V_A, f_B < f_A, t_B^s = t_A^s \frac{f_A}{f_B}$
 - For $N_B = \rho N_A$ samples: $E_B^s = \rho N_A(\hat{P}_A \frac{f_A}{f_B} t_A^s + C V_B^2 f_A \frac{f_A}{f_B} t_A^s)$
 - Reduced voltage:**
 - $V_C = \gamma V_A$, with $\gamma < 1$, $f_C = f_A, t_C^s = t_A^s$.
 - For $N_C = \rho N_A$ samples: $E_C^s = \rho N_A(\hat{P}_A t_A^s + C \gamma^2 V_A^2 f_A t_A^s)$.
- ρ is the oversampling needed to guard against faults.

Results

Energy ratios for cases B,C wrt to A showing saving regimes for different oversampling.



(1) K. Morris, F. Rizzi, B. Cook, P. Mycek, O. Le Maître, O. Knio, K. Sargsyan, K. Dahlgren and B. Debusschere, *Performance Scaling Variability and Energy Analysis for a Resilient ULFM-based PDE Solver*, Scal16, 2016.
 (2) K. Sargsyan, F. Rizzi, P. Mycek, O. Le Maître, O. Knio and B. Debusschere, *Fault resilient domain decomposition preconditioner for PDEs*, SIAM J. Sci. Comp., 37(5), 2015.
 (3) F. Rizzi, K. Morris, K. Sargsyan, P. Mycek, C. Saffa, O. Le Maître, O. Knio and B. Debusschere, *Exploring the Interplay of Resiliency and Energy Consumption for a Task-Based Partial Differential Equations Preconditioner*, PP4REE, co-located with PPoPP16, 2016.
 (4) P. Mycek, F. Rizzi, O. Le Maître, K. Sargsyan, K. Morris, C. Saffa, B. Debusschere and O. Knio, *Discrete a priori bounds for the detection of corrupted PDE solutions in exascale computation*, SIAM J. Sci. Comp., 39(1), 2017.