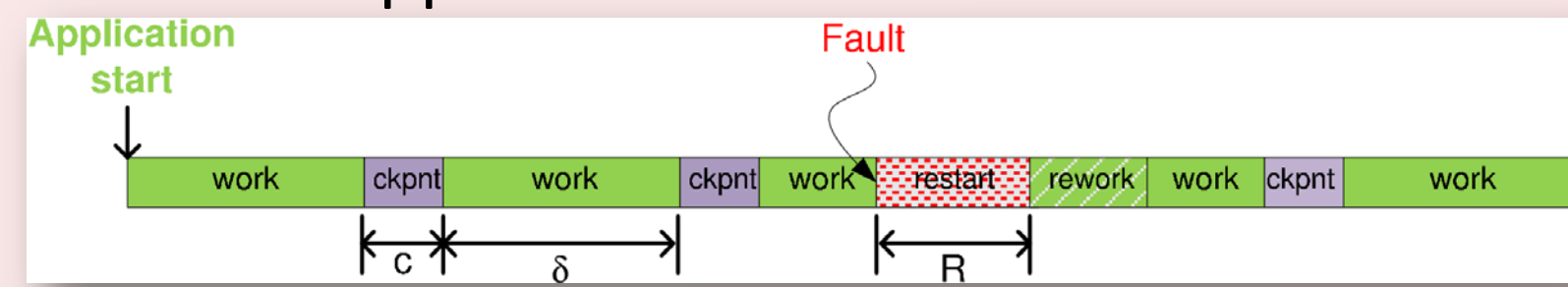


Investigating the Benefits of Redundancy Plus Checkpointing for Hard-Fault and Soft Error Protection in HPC

David Fiala, James Elliot, Kishor Kharbas Advisor: Frank Mueller (NCSTU)
 Collaborators: Christian Engelmann (ORNL), Rolf Riesen, Kurt Ferreira (SNL)

MOTIVATION

- Component failures require support of checkpoint/restart (C/R)
- Node or switch failure → application fails

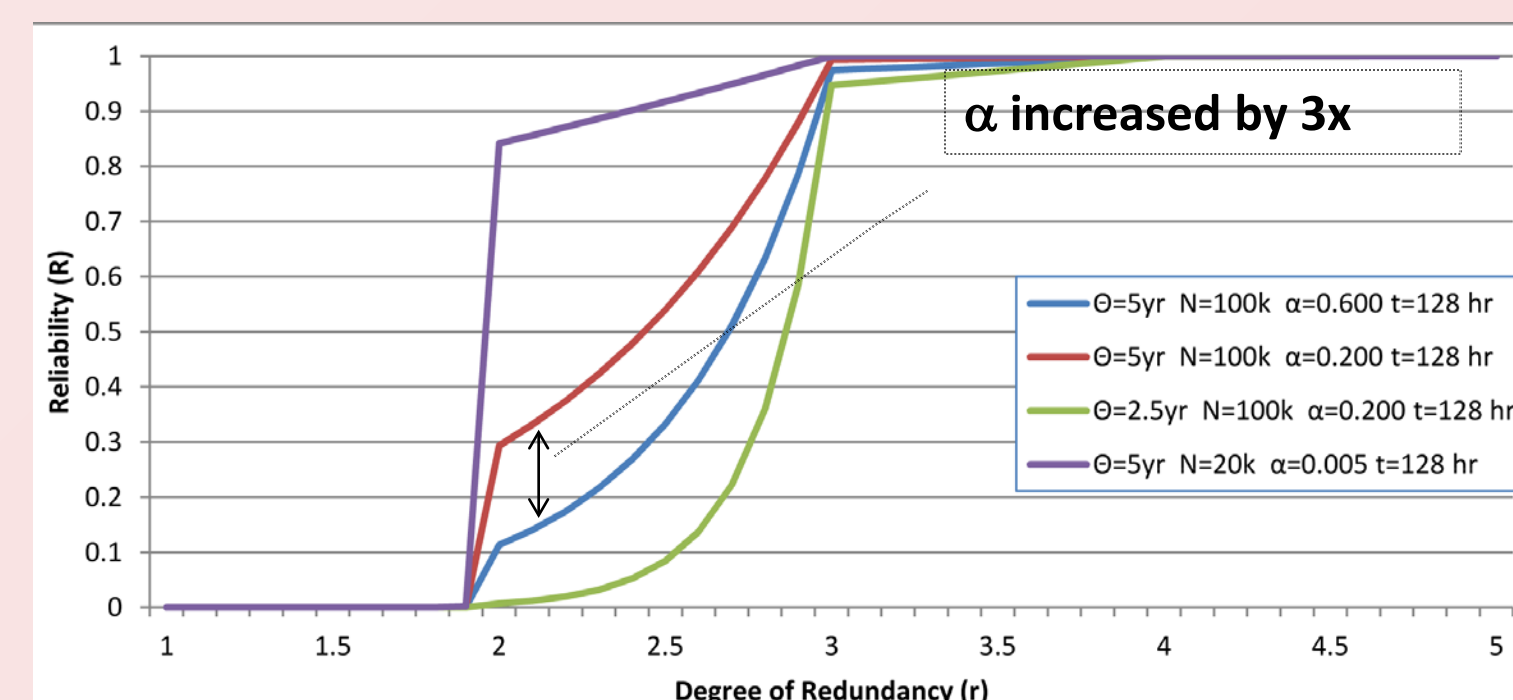


- Adding hardware increases the likelihood of faults
 - The probability of component failure combinatorially explodes
 - The mean-time-between-failure (MTBF) shortens
 - Overhead due to C/R increases exponentially
- Computation vs. overhead ratio can be between 85%-55%

168-hour Job, 5 year MTBF [Sandia]				
# nodes	Work	Checkpoint	Re-computation	Restart
100	96%	1%	3%	0%
1,000	92%	7%	1%	0%
10,000	75%	15%	6%	4%
100,000	35%	20%	10%	35%

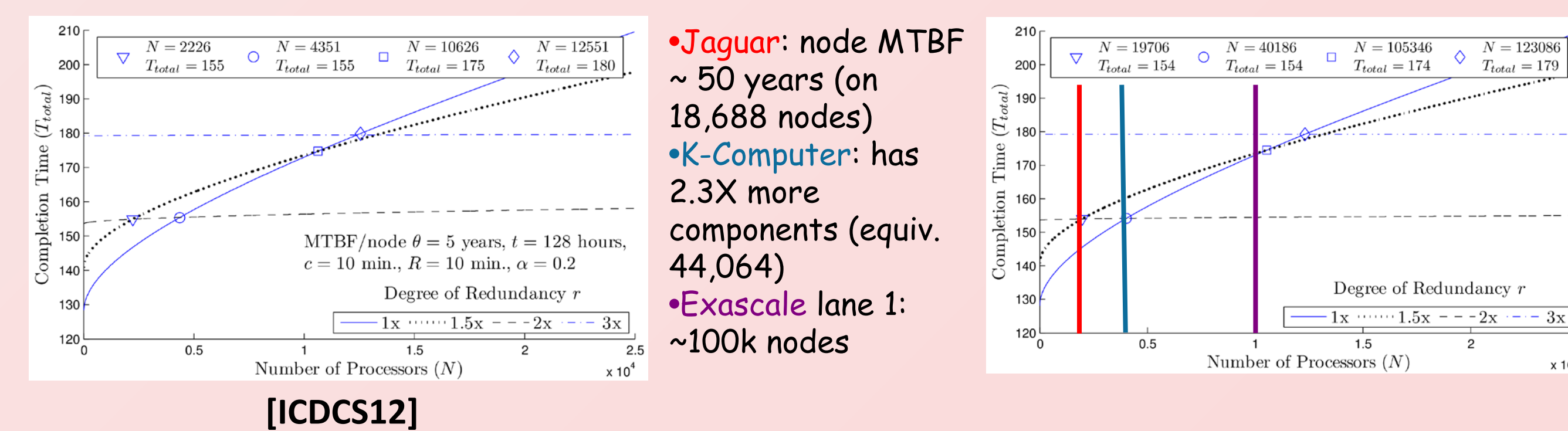
- C/R does not increase reliability, but redundancy can reverse this trend
 - Each redundant process decreases the probability of failure of replica processes
 - Less interruptions produces greater utilization
 - 100% redundancy provides 5x job throughput [Sandia]

- In the face of fail-stop failures, redundancy increases reliability



α = Mean time between failure (MTBF)
 α = Communication to computation ratio
 Time becomes function of α
 [ICDCS12]

- 2x redundancy may be beneficial now: at 78,536 processes, two dual redundant jobs of 128 hours can be run in the time of just one job without redundancy



- Another class of fault: **Silent Data Corruption (SDC)**
- SDC faults manifest themselves as bit-flips in storage or even within processing cores
 - In some cases bit-flips are not correctable or even detected
 - Exacerbating this situation, when SDC goes undetected invalid results are reported
 - Memory becomes corrupt, but applications continue to run
 - This is a severe problem for today's large-scale simulations

REDMPI: PROTECTING AGAINST SDC FAULTS

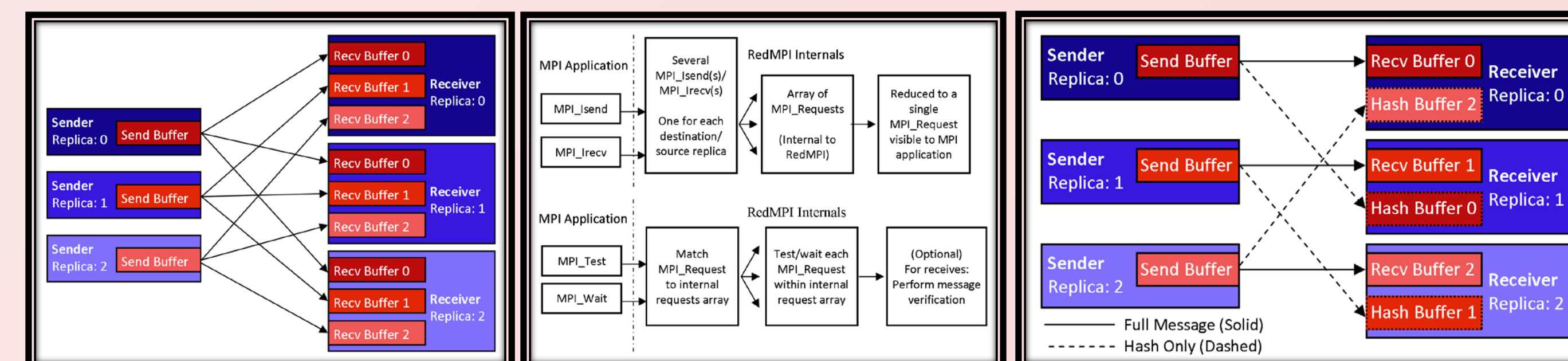
- Redundancy for HPC applications provides more than just fail-stop protection
 - The framework for providing redundancy also provides a platform to both protect applications and investigate the effects of faults

- Provide transparency by linking unmodified MPI applications with our RedMPI
- RedMPI provides redundancy to MPI applications by instrumenting the MPI profiling layer
 - Adjusted MPI rank and size provide illusion of normal rank numbers
 - SDC protection is afforded by augmenting MPI functionality to communicate with replicas

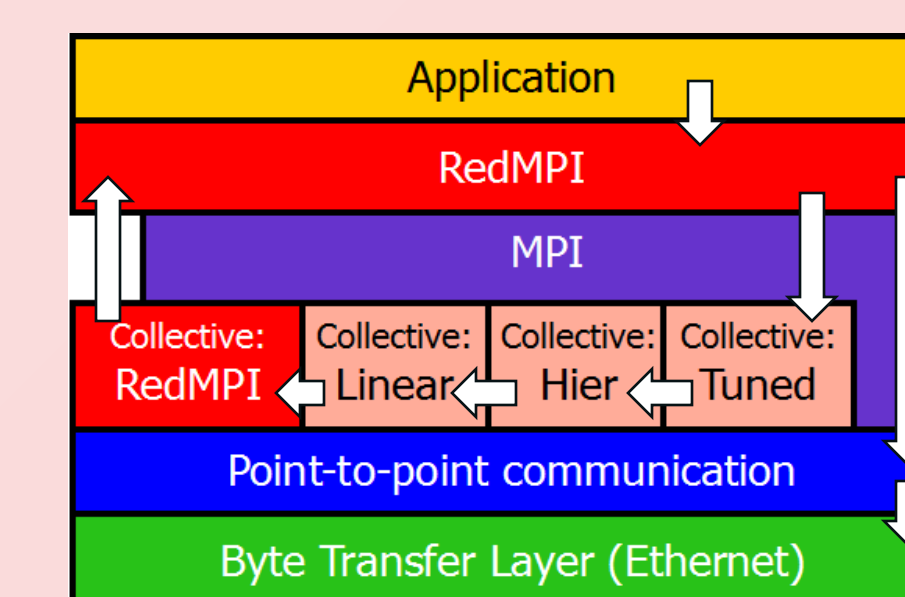


	No Redundancy	Dual Redundancy	Triple Redundancy or higher
Live SDC Detection	No	Yes	Yes
Live SDC Correction	No	No	Yes (via voting algorithm)

- Naïve SDC protection may be achieved by transmitting and comparing $r*r$ messages amongst r total replicas.
 - Induces high interconnect contention / bandwidth degradation
 - Compare received buffers, discard a mismatch



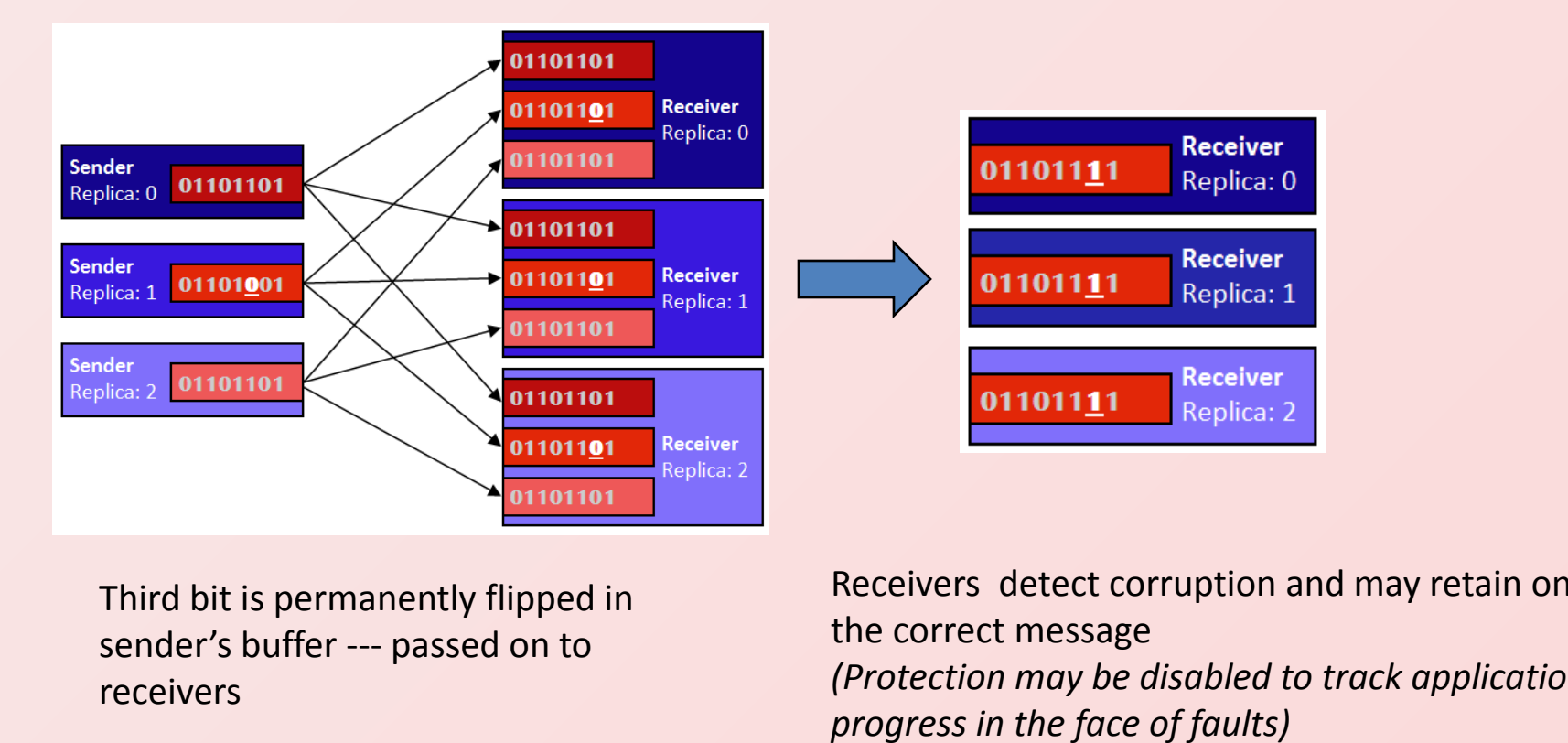
- Enhance performance by sending the original message plus a small hash to a separate replica
 - No longer dependent on $r*r$ communication
 - Comparison still performed on receiver-side
 - Hash mismatch triggers secondary voting protocol amongst receiving replicas



- Collectives require protection too – and must operate in the face of faults
 - Reuse existing topology-aware tuned algorithms, route messages via RedMPI to ensure success

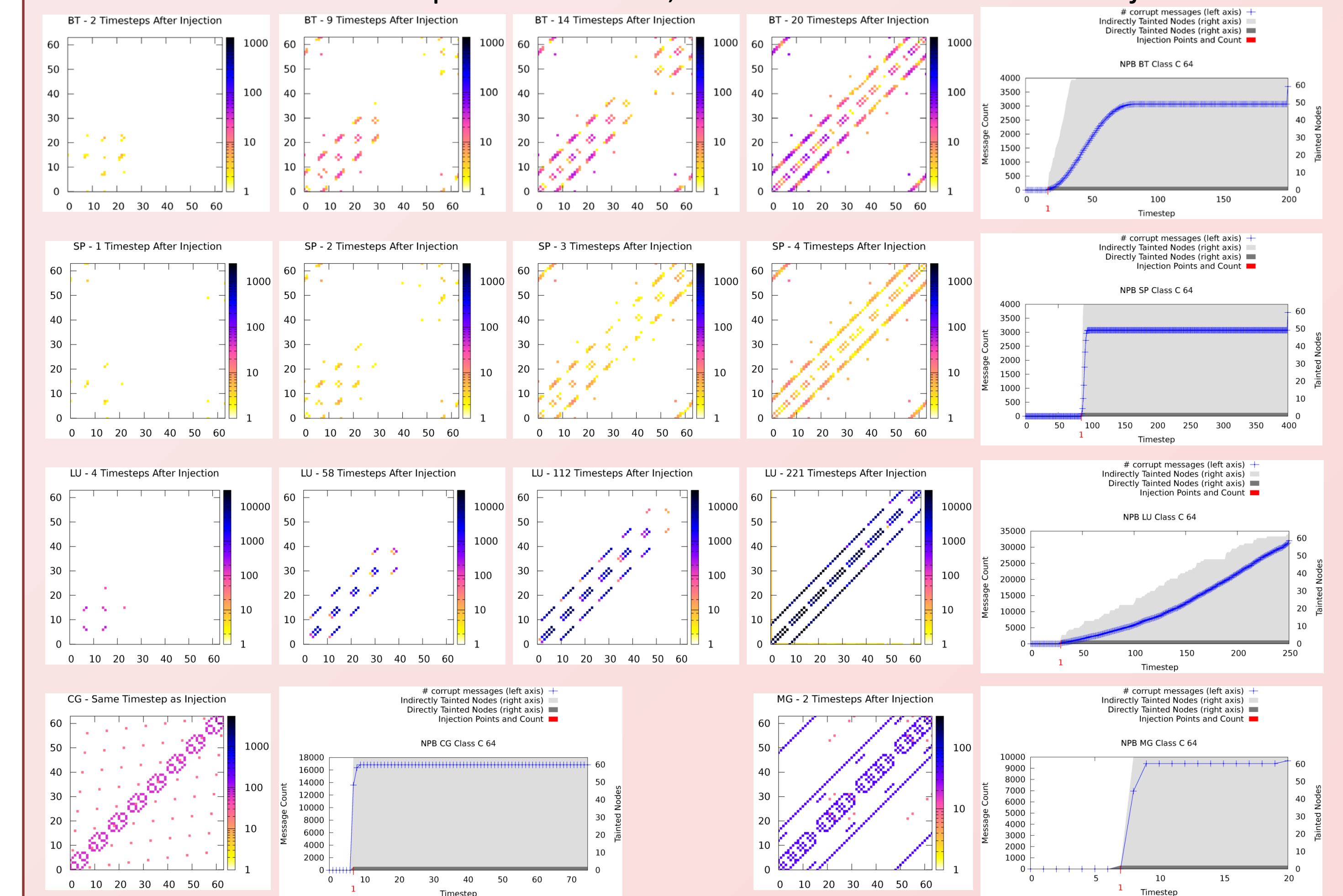
FAULT INJECTION

- Sender side: 1/x messages randomly receive single bit-flip in application buffer
 - This corrupts the sender's memory and spreads to peers as incorrect MPI messages



OBSERVATION: SDC PROPAGATION

- What is the impact of a just one bit-flip?
 - Run RedMPI with dual redundancy, but disable protection & track live execution
 - One set of correct replicas as control, one set tainted with an SDC injection



- Profound effects may spread quickly, depending on communication patterns

OBSERVED REDUNDANCY COSTS

- Experiments performed on 96 cluster nodes
 - AMD Opteron 6128 (Magny Core) – 16 cores per node – 32GB RAM per node
 - 40Gbit/s Infiniband for MPI Communication
- 1x: Uninstrumented Open MPI (No Redundancy)
- 2x: RedMPI with Dual Redundancy
- 3x: RedMPI with Triple Redundancy

LAMMPS - CHUTE-SCALED						NPB - CG					
Size	1x [sec]	2x [sec]	3x [sec]	2x OV	3x OV	Size	1x [sec]	2x [sec]	3x [sec]	2x OV	3x OV
128	137.5	138.4	139.0	0.6%	1.1%	128-D	201.42	205.87	215.51	2.2%	7.0%
256	138.3	140.4	140.0	1.6%	1.3%	256-D	127.21	132.61	136.64	4.2%	7.4%
512	139.2	140.2	141.0	0.7%	1.1%	512-D	70.10	77.54	83.67	10.6%	19.4%

SWEETSD						NPB - EP					
Size	1x [sec]	2x [sec]	3x [sec]	2x OV	3x OV	Size	1x [sec]	2x [sec]	3x [sec]	2x OV	3x OV
128	390.3	389.5	393.1	-0.2%	0.7%	128-D	72.31	72.63	72.74	0.4%	0.6%
256	428.2	427.5	431.2	-0.1%	0.7%	256-E	579.94	581.02	581.27	0.2%	0.2%
512	488.1	488.9	494.1	0.2%	1.2%	512-E	289.80	290.83	291.30	0.4%	0.5%

HPCCG - wildcard receives present						NPB - FT					
Size	1x [sec]	2x [sec]	3x [sec]	2x OV	3x OV	Size	1x [sec]	2x [sec]	3x [sec]	2x OV	3x OV
128	99.8	99.8	125.8	0.0%	26.0%	128-C	117.45	117.95	118.68	0.43%	1.05%
256	99.6	128.8	131.0	29.3%	31.5%	256-C	68.82	68.62	71.77	-0.29%	4.29%
512	126.4	146.2	152.3	15.7%	20.5%	512-D	222.75	228.76	234.97	2.70%	5.49%

CONCLUSIONS

- For large systems, C/R + redundancy increases job throughput
- Redundancy is cheap in terms of software overhead
- Application sensitivity to soft errors may be very high
- SDC protection comes free as redundancy is used to increase system resilience
 - RedMPI can successfully protect applications from SDC faults and continue execution to a successful, correct completion