

Fenix: An Online Failure Recovery Library for MPI applications on top of ULFM

**Marc Gamell¹, Keita Teranishi²,
Rob Van Der Wijngaart³, Manish Parashar¹**

¹ Rutgers Discovery Informatics Institute (RDI²), Rutgers University

² Sandia National Laboratories

³ Intel Corporation

Based on:

“Exploring Automatic, Online Failure Recovery for Scientific Applications at Extreme Scales”, SC14
Marc Gamell, Daniel S. Katz, Hemanth Kolla, Jacqueline Chen, Scott Klasky, Manish Parashar

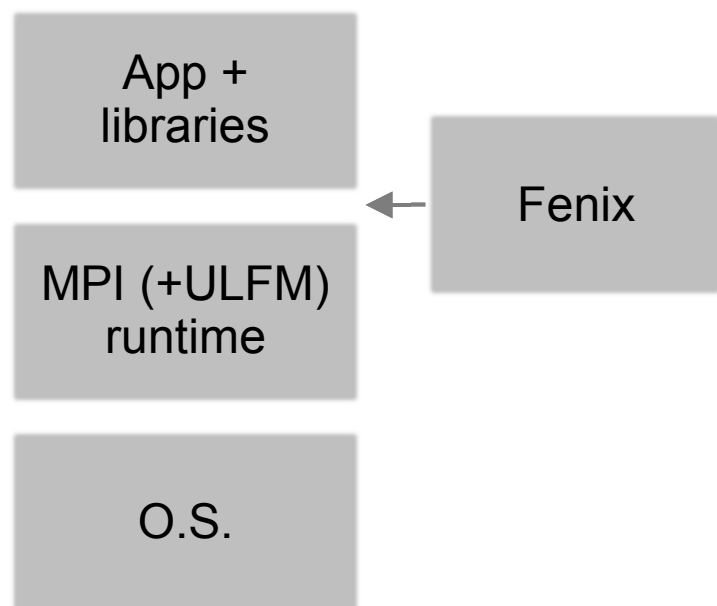
Introduction / Goal

- Fenix can be used effectively, efficiently, and productively to provide online fault tolerance to a major production code (S3D)
- How can an application use Fenix?
- How does Fenix use ULFM and PMPI?

Table of Contents

- Introduction
- Key Contributions
- Motivating Use Case
- Process Recovery Interface
- Data Recovery Interface
- Usage Example
- Experimental Evaluation
- Conclusion

Key Contributions



Implementation details

- Built on top of **ULFM**
- Tested up to
 - 16384 cores w/ failures
 - 250k cores w/o failures
- Available for C and Fortran applications

Approach

- Targets **MPI-based** parallel applications
- Offers two disjoint interfaces

Fenix Interfaces:

1. Process/rank recovery

- **Online, semi-transparent recovery from process, node, blade and cabinet failures**
- **Tolerates a variety of MTBFs**
 - even extreme MTBFs of <1 minute

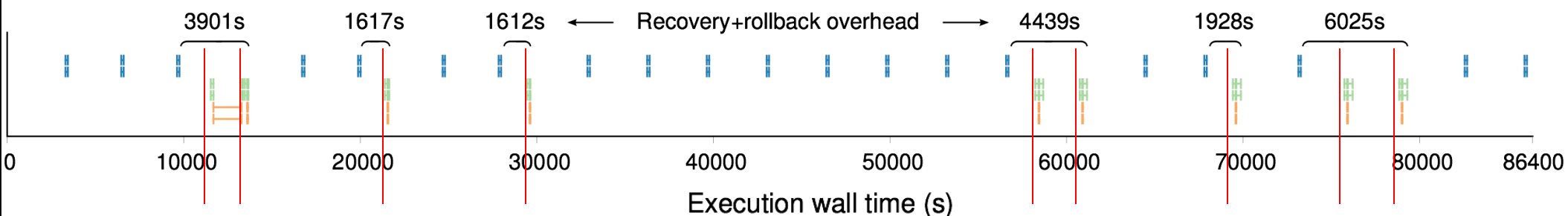
2. Data recovery

- Uses **application-specific, double in-memory, implicitly coordinated checkpoints**

Experimental Evaluation

- Deployed Fenix on **Titan** Cray XK7 at ORNL
- **S3D** combustion numerical simulation
- Sustained performance with **MTBF = 47 seconds**
- Experiments inject real process failures (SIGKILL)

Motivating Use Case – S3D production runs



- 24-hour tests using Titan (125k cores)
- **9 process/node failures** over 24 hours
- Failures are promoted to **job failures**, causing all 125k processes to exit
- Checkpoint (5.2 MB/core) has to be done to the PFS

Checkpoint (per timestep)

Restarting processes

Loading checkpoint

Rollback overhead

Total overhead

Total cost

55 s 1.72 %

470 s 5.67 %

44 s 1.38 %

1654 s 22.63 %

31.40 %

Motivating Use Case – Possible solution

- Process failures cannot be promoted to job failures, to:
 - **Reduce recovery** cost
 - Keep process memory (contains checkpoints)

Online recovery

- **Checkpoint frequency** has to be dramatically **increased**, e.g.

**Checkpoint application-specific data
in process memory w/o coordination**

ULFM – User Level Failure Mitigation

User Level Failure Mitigation is a **set of MPI extensions to report errors, provide interfaces to stabilize the distributed state, and restore the communication capabilities in applications affected by process failures**. Relevant communicators, RMA windows and I/O files can be reconstructed online, without restarting the application, as required by the user recovery strategy.

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.

Two Philosophy Changes for the Application

Fenix:

- Provides semi-transparent online recovery to particular application types
- Implements a particular recovery mode thanks to generic ideas behind ULFM

Application view:

- Offer a **single point** in which all ranks will **return** upon failure has been recovered
 - longjump
- Application should use **Fenix's resilient communicator** instead of MPI_COMM_WORLD
 - Resilient communicator is created when initializing Fenix
 - All communicators derived from it can also be used

What about other libraries?

- If they do not use Fenix (and keep state): teardown and re-initialize
- If they use Fenix themselves:
 - Not supported for the moment
 - Maybe libraries can register callbacks to be called upon recovery

Fenix – Failure Recovery Interface

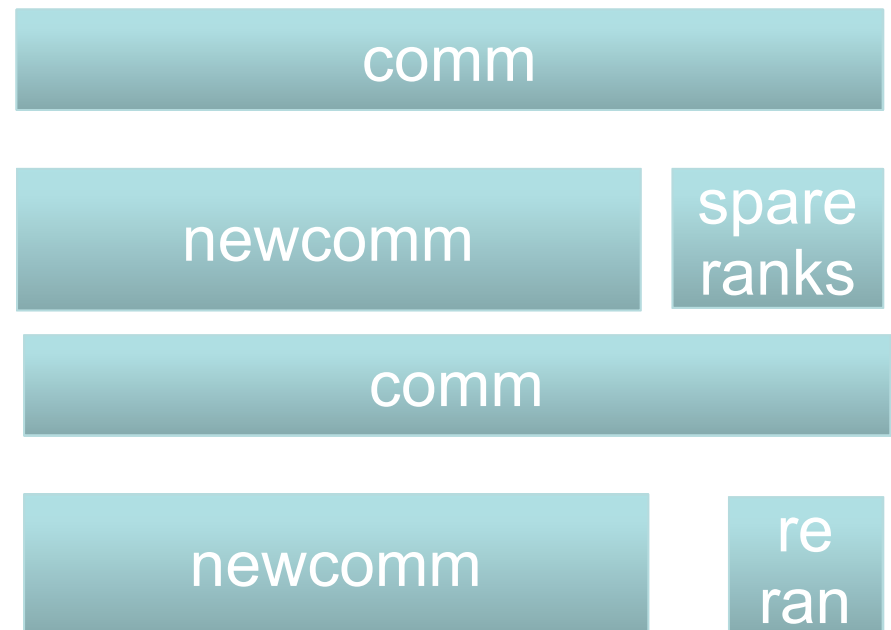
```

void Fenix_Init ( int *status,
MPI_COMM_WORLD → MPI_Comm comm,
MPI_Comm *newcomm,
App should use ← MPI_Comm *newcomm,
newcomm instead of MPI_COMM_WORLD
int *argc, int ***argv,
int num_spare_ranks,
Fenix_Comm_repair_policy repair_policy, ← NO_SPAWN
int *error);                               ← SPAWN
    
```

```
void Fenix_Finalize ( );
```

newcomm: resilient communicator

- Non-shrinking: SPAWN policy
- Shrinking: NO_SPAWN + num_spare_ranks = 0
- Mixed (only spares): NO_SPAWN policy + num_spare_ranks > 0



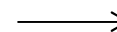
Fenix + ULFM – Recovery Stages

1. Failure detection

- With ULFM, MPI communicating calls may return failure codes
- Fenix captures them using MPI profiling interface
 - ☹ For now, this implies **no tools can be used in conjunction with Fenix**
 - ☺ In the future, maybe **MPI will replace PMPI with a method that allows attaching multiple tools (QMPI/MPI Extension Interface)?**
- Uses MPI_Comm_revoked to spread notification

2. Environment recovery

- Repair only main communicator
 - Non-shrinking model:
 - Use spare process pool and
 - Re-spawn processes
 - **Non-shrinking up to a certain number of failures**
 - **Use spare process pool**
 - Shrinking model
- Re-creation of user communicators: done by the user



This is shown in the following slides

Fenix + ULFM – Internal Initialization Pseudocode

Assume we have two communicators:

- MPI_Comm comm (originally, a dup from Fenix_Init's input comm)
- MPI_Comm newcomm (repaired communicator to be used, returned as newcomm from Fenix_Init)

```
fenix_internal_init() {  
    MPI_Comm_set_errhandler(comm, MPI_ERRORS_RETURN);  
    MPI_Comm_rank(comm, &rank);  
    MPI_Comm_rank(comm, &size);  
  
    if(rank >= size - num_of_spare_ranks) { // SPARE RANK  
        PMPI_Comm_split(comm, MPI_UNDEFINED, rank, &newcomm);  
  
        for(;;) { // wait for failure to occur  
            ret = MPI_Recv(..., MPI_ANY_SOURCE, MPI_ANY_TAG, comm);  
            if(ret == MPI_SUCCESS)  
                exit();  
        }  
  
    } else { // NOT A SPARE RANK  
        PMPI_Comm_split(comm, 0, rank, &newcomm);  
    }  
}
```

Fenix + ULFM – Failure Detection Pseudocode

- Override MPI functions and use PMPI interface, e.g.

```
int MPI_Send(void *buf, int count, MPI_Datatype datatype,
             int dest, int tag, MPI_Comm comm)
{
    int ret;
    ret = PMPI_Send(buf, count, datatype, dest, tag, comm);
    fenix_internal_test_return_code(ret);
    return ret;
}
```

Fenix + ULFM – Testing for Failures Pseudocode

```
void fenix_internal_test_return_code(int ret) {  
    switch(ret) {  
        case MPI_SUCCESS:  
            return;  
        case MPIX_ERR_PROC_FAILED:  
            MPIX_Comm_revoke(comm);  
            MPIX_Comm_revoke(newcomm);  
            fenix_internal_revoke_user_communicators();  
        case MPIX_ERR_REVOKED:  
            fenix_internal_communicator_repair();  
            break;  
    }  
}
```

Fenix + ULFM – Recovery Pseudocode (1/2)

```
fenix_internal_communicator_repair() {  
    MPI_Comm comm_shrink;  
    MPIX_Comm_shrink(comm, &comm_shrink);  
  
    // Do we have enough spare ranks?  
    MPI_Comm_size(comm, &old_comm_size);  
    MPI_Comm_size(comm_shrink, &new_comm_size);  
    ranks_needed = old_comm_size - new_comm_size;  
    if(num_of_spare_ranks < ranks_needed)  
        // for today's discussion, consider an error  
  
    // Which ranks failed?  
    MPI_Comm_rank(comm, &old_rank);  
    PMPI_Allgather(&old_rank, 1, MPI_INT,  
                  survivor_ranks, 1, MPI_INT, comm_shrink);  
    // locally determine the failed ranks by using the survivor_ranks array  
  
    ...  
}
```

Fenix + ULFM – Recovery Pseudocode (2/2)

```
fenix_internal_communicator_repair() {  
    // Shrink the communicator  
    // Do we have enough spare ranks?  
    // Which ranks failed?  
    ...  
  
    // Assign spare ranks to failed ranks  
    if(old_rank >= size_newcomm) // WAS A SPARE RANK  
        if((old_comm_size-1-old_rank) < procs_needed) // RECOVERED RANK  
            old_rank = // pick one of the failed ranks  
    else // WAS NOT A SPARE RANK  
        PMPI_Comm_free(&newcomm);  
    num_of_spare_ranks -= num_failed_ranks;  
  
    // Re-name the spare ranks in the new communicator  
    PMPI_Comm_free(&comm);  
    PMPI_Comm_split(comm_shrink, 0, old_rank, &comm);  
    PMPI_Comm_free(&comm_shrink);  
  
    fenix_internal_init();  
  
    // no matter who called fenix_internal_communicator_repair, return to Fenix_Init()  
}
```

Which communicators does Fenix revoke?

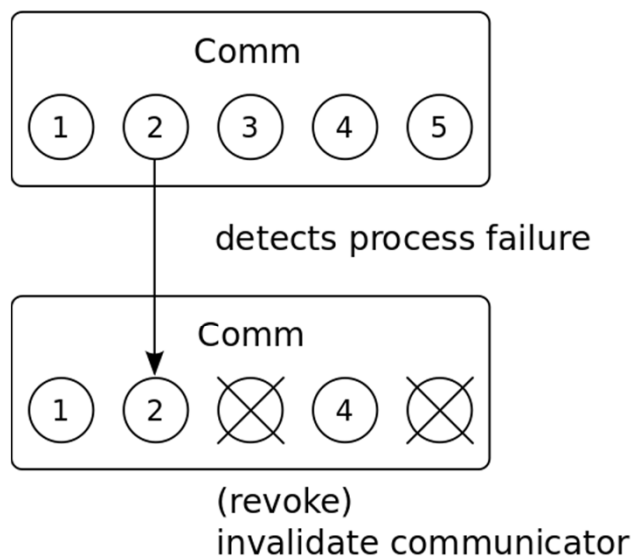
- When failure is detected by a rank, ALL communicators derived from newcomm are revoked
 - The only communicators that the application should use with Fenix philosophy
- Derived communicators are detected by the profiling interface
 - MPI_Comm_split, MPI_Comm_create, MPI_Comm_dup...
 - Captured by profiling interface when they are first derived
 - Upon failure, revoked

Fenix – Recovery Stages

1. Failure detection

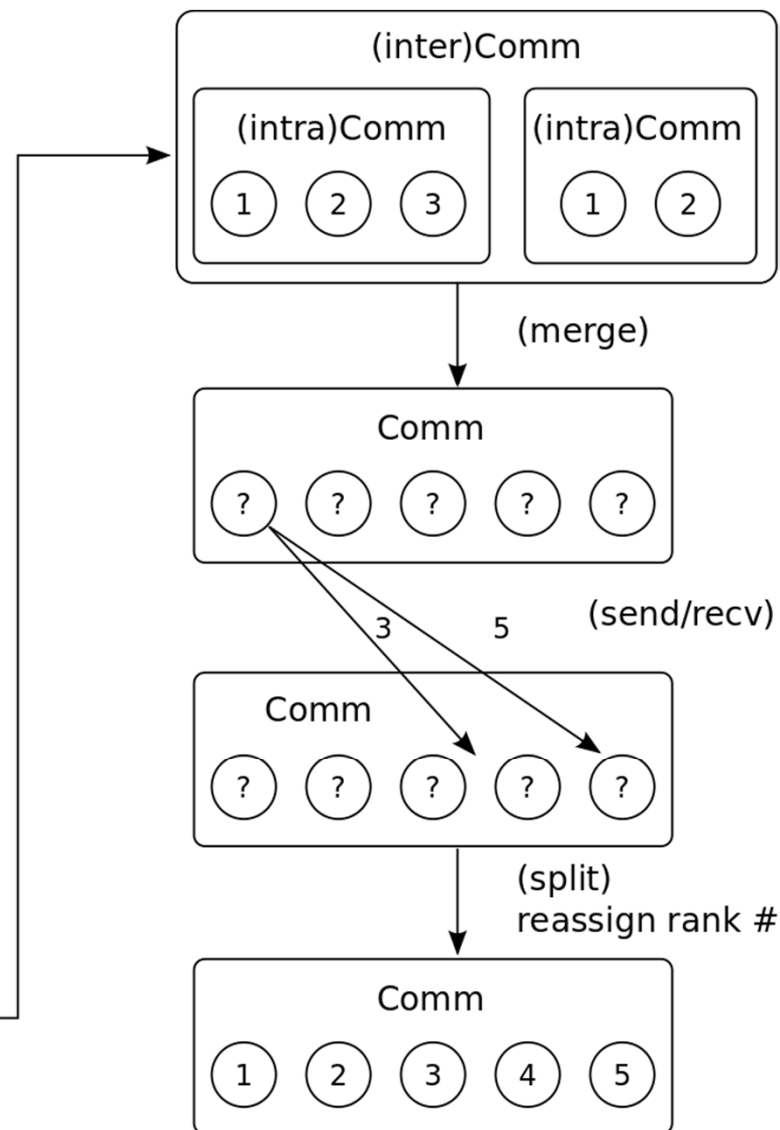
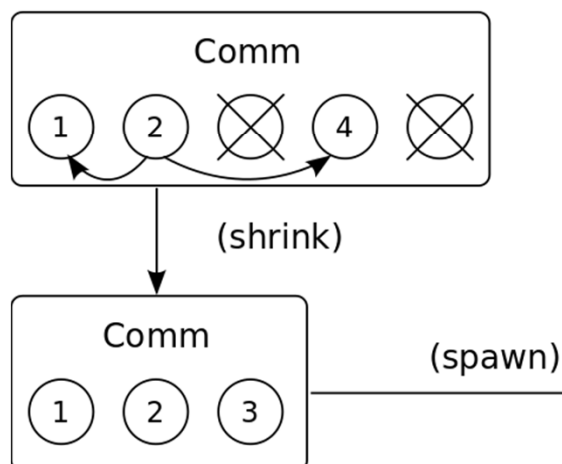
Based on

- ULFM return codes
- MPI profiling interface
- Comm revoke



2. Environment recovery

- Re-spawn / Spare process pool
- Repair “newcomm” communicator
- Delay re-creation of user comms



“Respawn” recovery mode depicted (similar for spare process pool)

Table of Contents

- Introduction
- Key Contributions
- Motivating Use Case
- Process Recovery Interface
- **Data Recovery Interface**
- Usage Example
- Experimental Evaluation
- Conclusion

Core of Fenix Data Recovery Interface

```
Fenix_Checkpoint ( int member_id,  
                  void *buffer,  
                  int count,  
                  MPI_Datatype datatype,  
                  Fenix_Checkpoint_group group,  
                  Fenix_Checkpoint_subset subset);
```

```
Fenix_Checkpoint_commit (  
    int *index,  
    Fenix_Checkpoint_group group);
```

Table of Contents

- Introduction
- Key Contributions
- Motivating Use Case
- Process Recovery Interface
- Data Recovery Interface
- **Usage Example**
- Experimental Evaluation
- Conclusion

Fenix – Simple Usage Example

Assume work1() and work2()
use an MPI communicator to
communicate with other ranks

```
1  /* Non-fault-tolerant version */
2  int main()
3  {
4      int it;
5      int A[100], B[50];
6
7      initialize(A, B);
8
9      for(it=0 ; it<1000 ; it++) {
10         work1(A, MPI_COMM_WORLD);
11         if(A[0] > 200) {
12             work2(A, B, MPI_COMM_WORLD);
13         }
14     }
15 }
```

Rank failure line 34
(due to imbalance, some ranks
are in it=102, the rest are in
it=103 or it=104)

Fenix will resume the execution,
returning from Fenix_Init()

```
1  /* Fault tolerant version with Fenix */
2  int main()
3  {
4      int it;
5      int A[100], B[50];
6      int status;
7      MPI_Comm new_comm_world;
8
9      Fenix_Init(&status, MPI_COMM_WORLD, &new_comm_world,
10         ...,
11         10, // num_spare_ranks
12         FENIX_COMM_REPAIR_POLICY_NO_SPAWN, // repair_policy
13         ...);
14     if( status == FENIX_STATUS_INITIAL_RANK ) {
15         /* no failure occurred */
16         it = 0;
17         initialize(A, B);
18         Fenix_Checkpoint(990, &it, 1, MPI_INT);
19         Fenix_Checkpoint(991, A, 100, MPI_INT);
20         Fenix_Checkpoint(992, B, 50, MPI_INT);
21         Fenix_Checkpoint_commit();
22     } else {
23         /* ranks recovered from a failure, now restore data */
24         Fenix_Checkpoint_restore(990, &it, 1, MPI_INT,
25             FENIX_LATEST_COMMIT);
26         Fenix_Checkpoint_restore(991, A, 100, MPI_INT,
27             FENIX_LATEST_COMMIT);
28         Fenix_Checkpoint_restore(992, B, 50, MPI_INT,
29             FENIX_LATEST_COMMIT);
30
31         for( ; it<1000 ; it++) {
32             Fenix_Checkpoint(990);
33             work1(A, new_comm_world);
34             if(A[0] > 200) {
35                 work2(A, B, new_comm_world);
36                 Fenix_Checkpoint(992);
37             }
38             Fenix_Checkpoint(991);
39             Fenix_Checkpoint_commit();
40         }
41     }
42 }
```

This sets it=101,
and A, B have the
same values they
had at iteration 101

VS

S3D Modifications

Main function

Only 35 new,
changed, or
rearranged lines
in S3D code

Topology module

VS

Fenix-enabled

Table of Contents

- Introduction
- Key Contributions
- Motivating Use Case
- Process Recovery Interface
- Data Recovery Interface
- Usage Example
- **Experimental Evaluation**
- Conclusion

Goal and Methodology

- Experiment:

S3D execution mimicking a future extreme-scale scenario

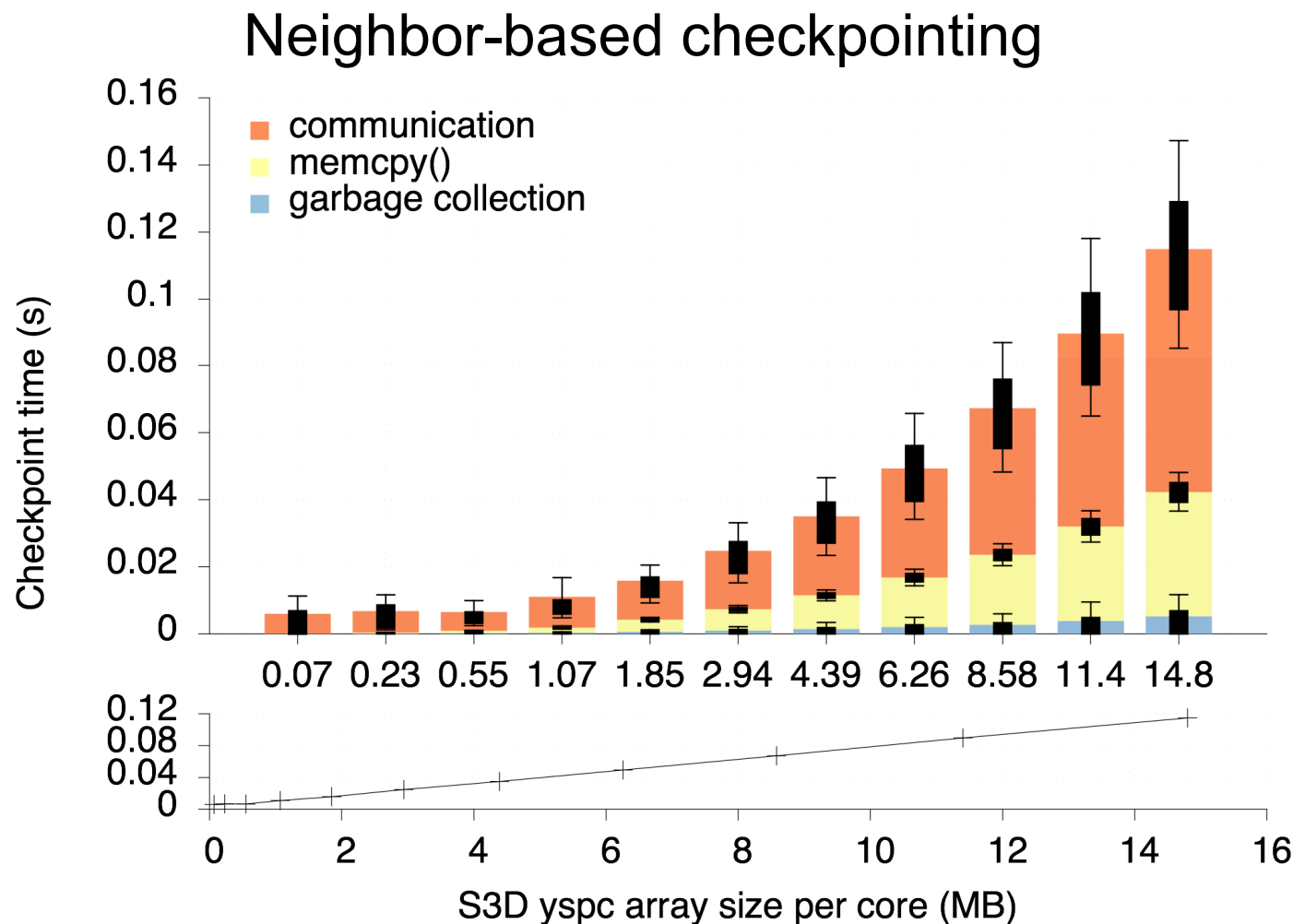
- Injecting failures every (MTBF):
 - 47 seconds
 - 94 seconds
 - 189 seconds
- Checkpoint 8.58 MB/core

- Procedure:

1. Evaluate checkpoint scalability
2. Calculate and validate optimal checkpoint interval
3. Evaluate recovery scalability
4. Run experiment

- Evaluation on ORNL Titan - Cray XK7

1. Failure-free checkpoint cost (data size)

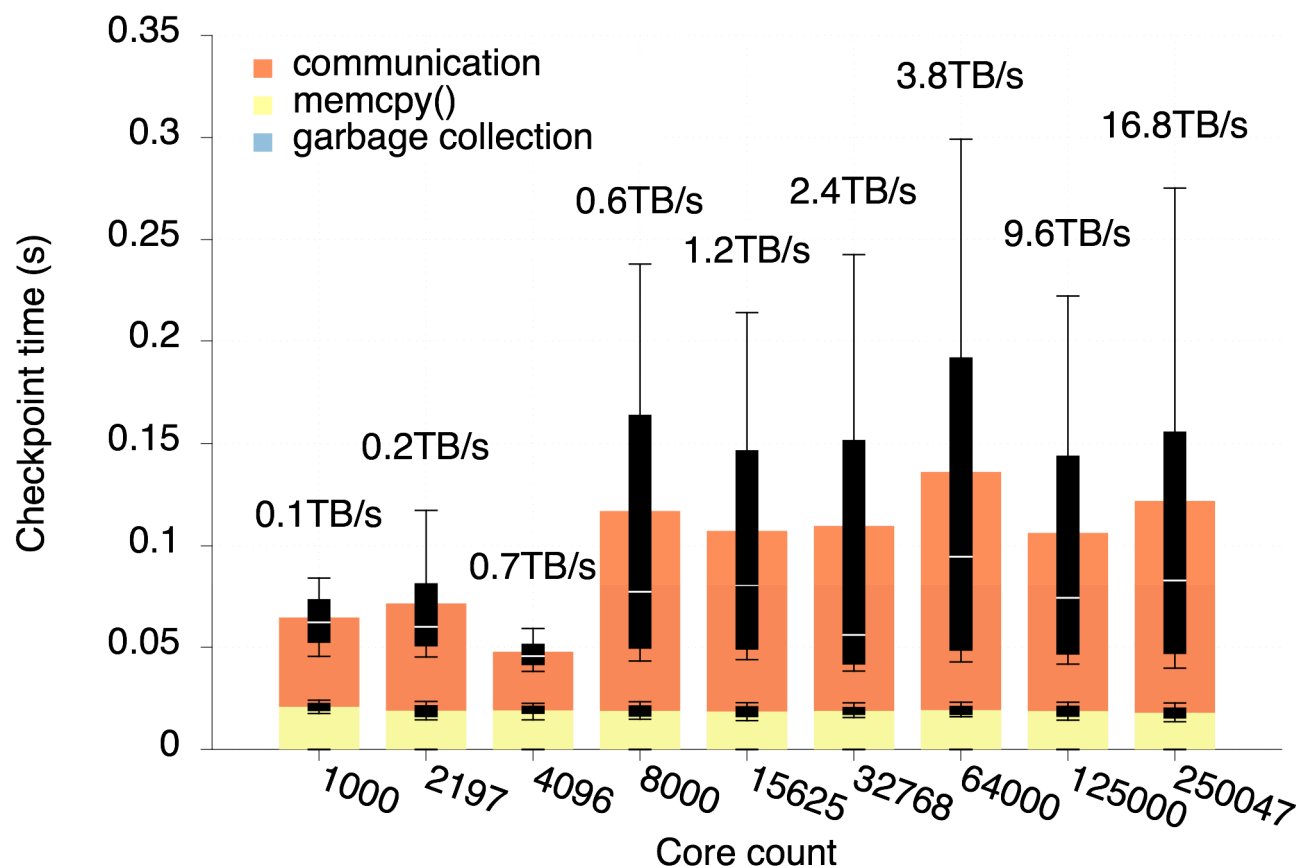


Conclusions:

- Scales linearly with data size increase
- Huge communication cost

1. Failure-free checkpoint cost (core count)

Neighbor-based checkpointing

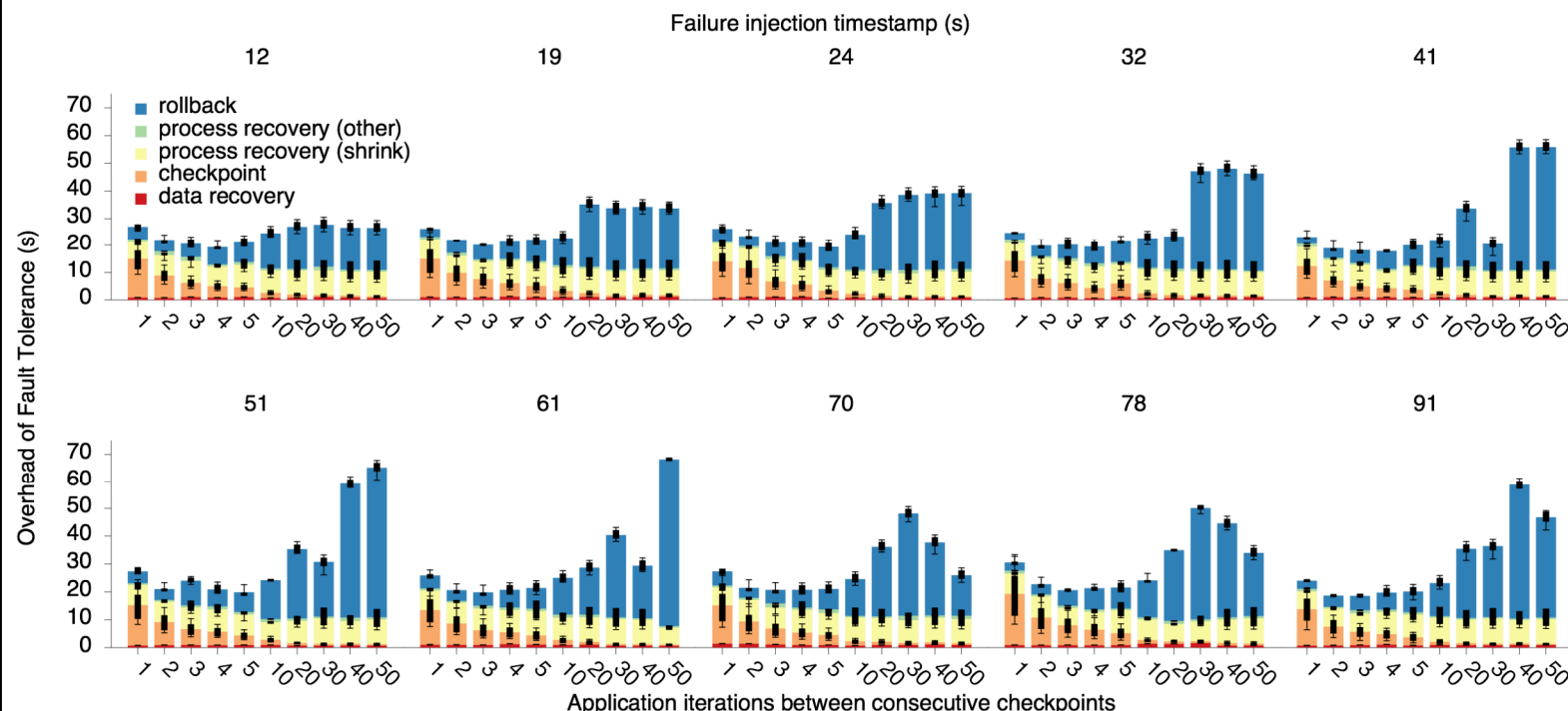


Conclusions:

- **Good scalability**
- **Really small footprint: O(0.1s)**

2. Optimal **checkpoint rate**

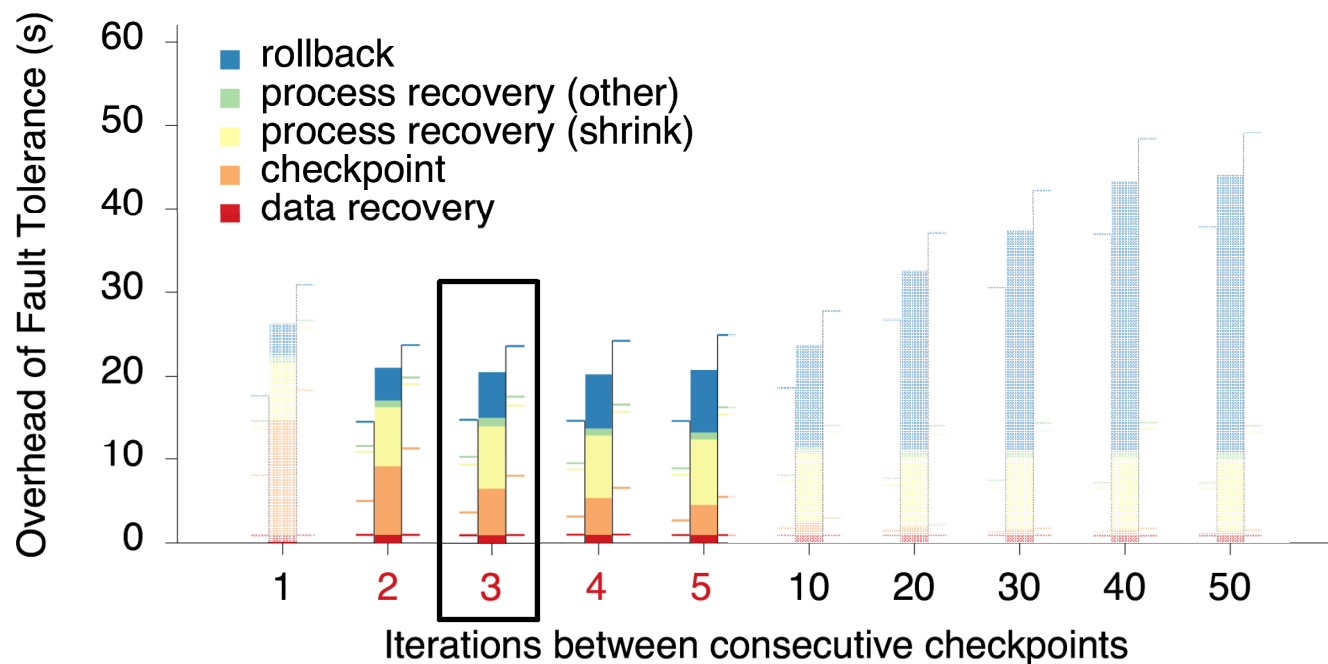
- Calculated by Young's formula:
 - T_S = checkpoint time = 0.0748 s (with 2197 cores)
 - $T_C \equiv \sqrt{\frac{\text{system's MTBF}}{2 I_S I_F}} = \{47, 94, 189\} \text{ s}$ (3 tests)
 - = {2, 3, 4} S3D iterations
- Empirically validating the 94-second MTBF optimal rate:
 - Inject only one failure at a specific wall time within the first 94 seconds of a test
 - Repeat using:
 - Different checkpoint rates
 - Different failure injection timestamps
 - Optimal rate must offer smaller overall rollback cost!



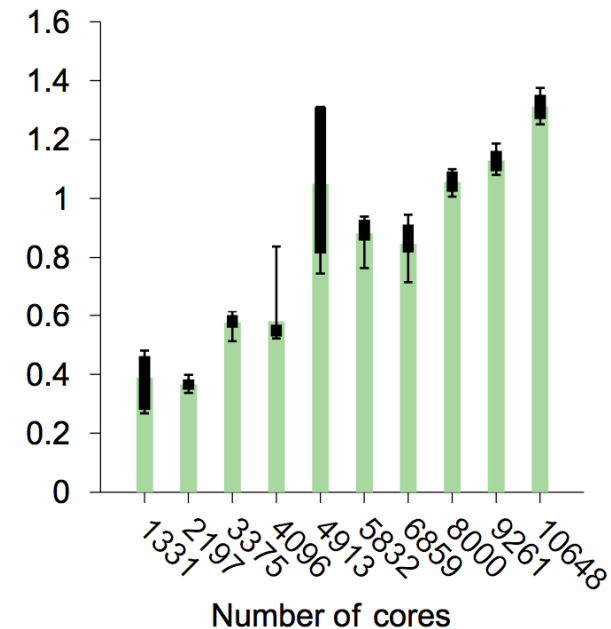
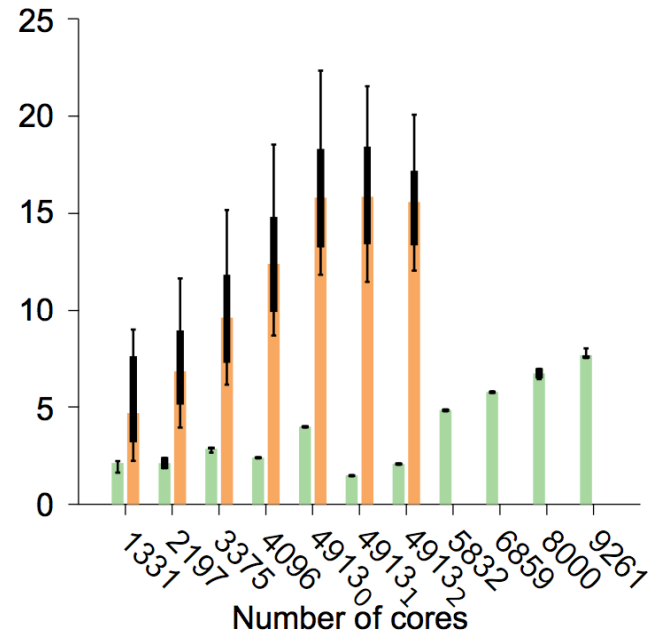
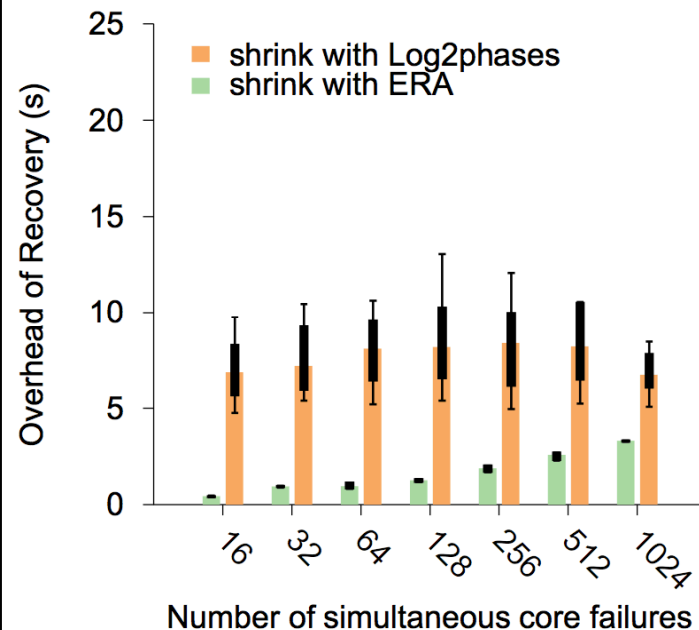
Average of all failure injection timestamps



2 to 5 iterations is the optimal checkpoint period, validating Young's formula



3. Recovery overhead

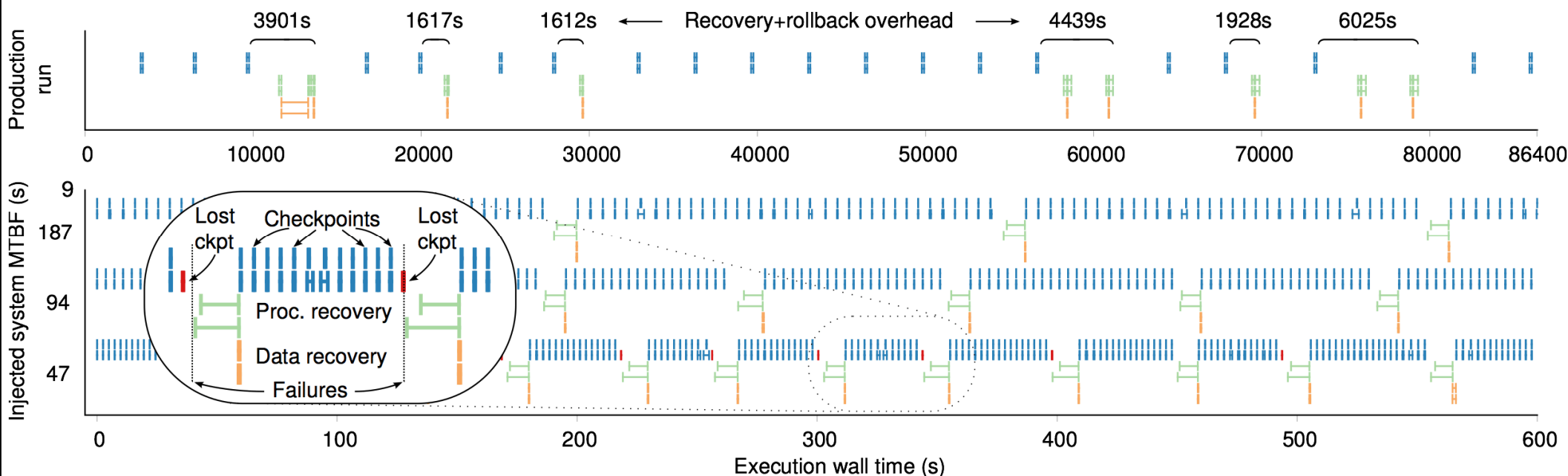


- 2197 cores

- 256-core failure
(i.e. 16 nodes)

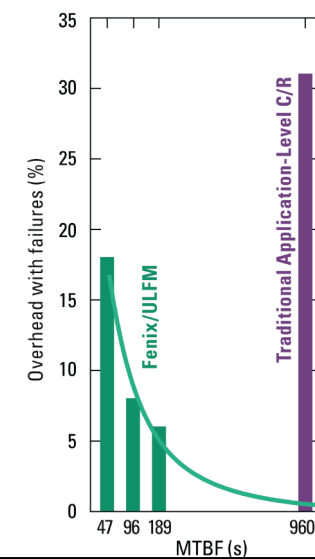
- 16-core failure
(i.e. 1 nodes)

4. Recovering from high-frequency failures

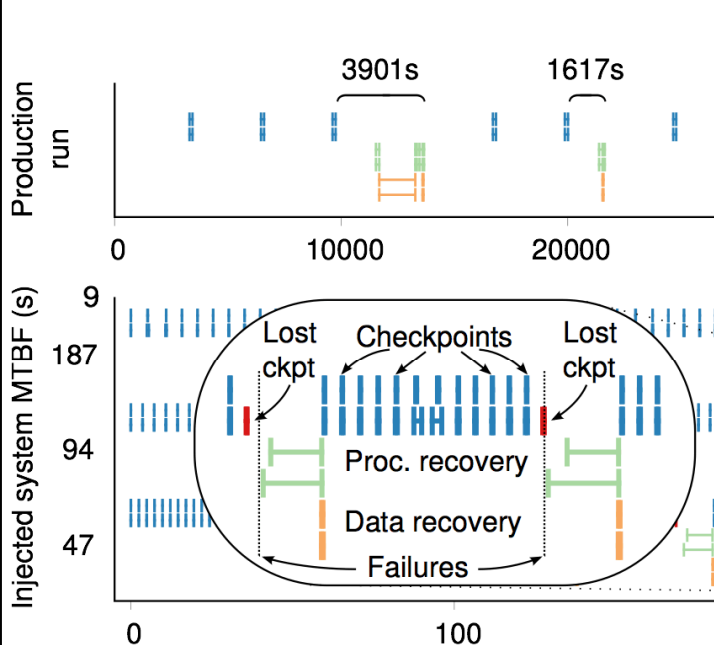


Conclusions:

- Online recovery allows the usage of in-memory checkpointing, $O(0.1s)$
- Efficient recovery from high-frequency node failures, as exascale compels



4. Recovering from high-frequency failures



Conclusions:

- Online recovery allows memory checkpointing
- Efficient recovery from failures, as exascale c

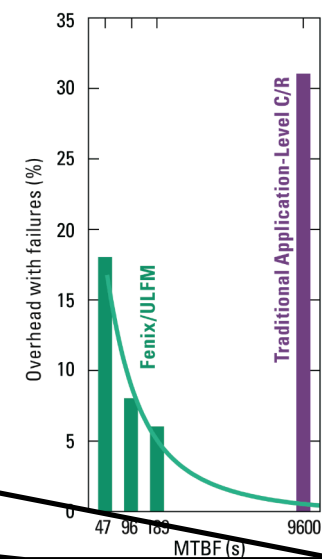
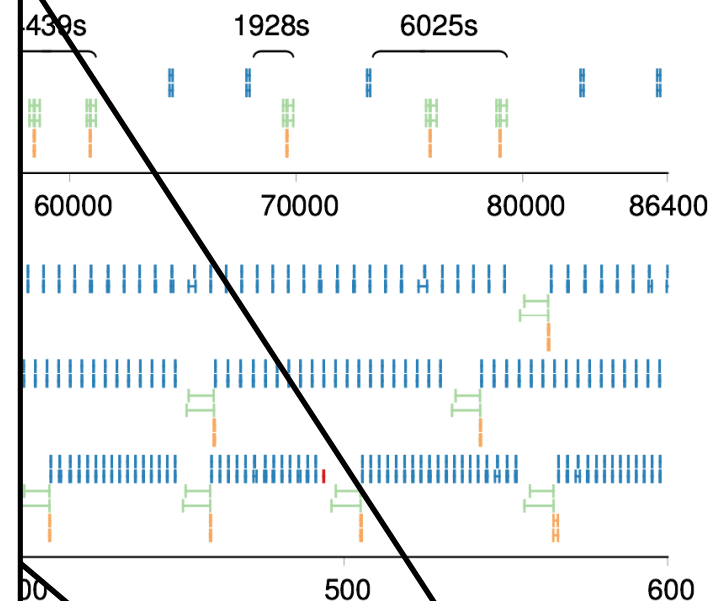
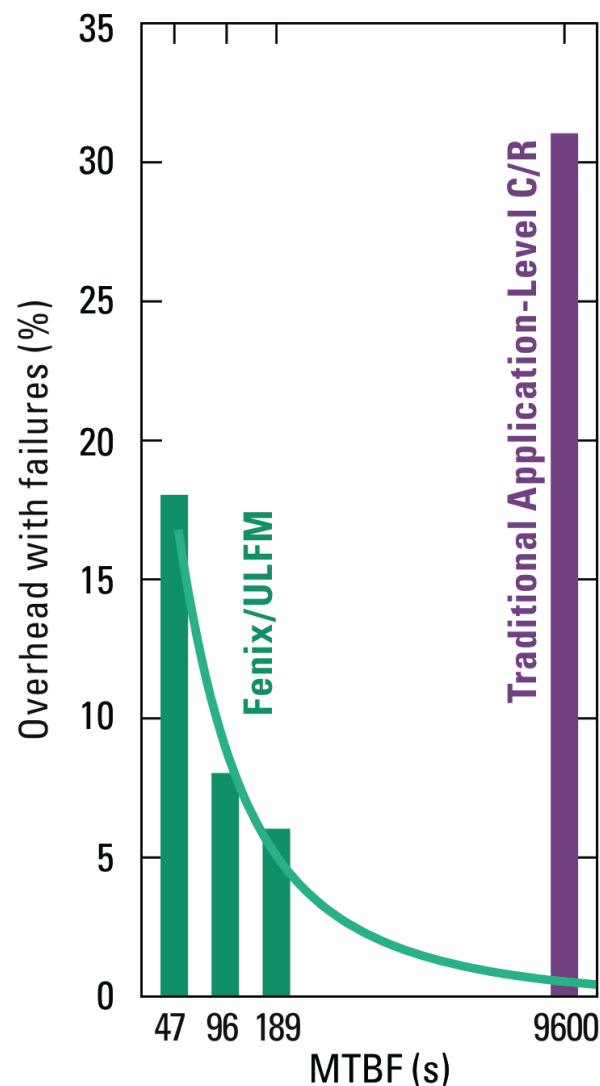


Table of Contents

- Introduction
- Key Contributions
- Motivating Use Case
- Process Recovery Interface
- Data Recovery Interface
- Usage Example
- Experimental Evaluation
- **Conclusion**

Conclusion

- Application-awareness can help resilience at scale
- Fenix provides
 - **Online Failure Recovery (reducing failure overhead)**
 - (uses in-memory, application-specific, **high frequency checkpointing**)
 - **Simple API built on top of ULFM:** only 35 new, changed, or rearranged lines in S3D code
- Deployed and empirically tested on Titan Cray XK7
 - **S3D+Fenix tolerate failure rates <1 min** with lower overhead as coordinated CR with failure rates of ~2.5 hours

Thank you