

ULFM

Обзор функциональных возможностей

User Level Failure Mitigation

The User Level Failure Mitigation (ULFM) proposal was developed by the MPI Forum's Fault Tolerance Working Group [1].

Designing the mechanism that users would use to manage failures was built around three concepts:

1. **simplicity**, the API should be easy to understand and use in most common scenarios;
2. **flexibility**, the API should allow varied fault tolerant models to be built as external libraries and;
3. **absence of deadlock**, no MPI call (point-to-point or collective) can block indefinitely after a failure, but must either succeed or raise an MPI error.

[1] <http://fault-tolerance.org/>

User Level Failure Mitigation

To use this ULFM implementation, an MPI application must change the default error handler on (at least)

❑ **MPI_COMM_WORLD** from **MPI_ERRORS_ARE_FATAL**

to either

❑ **MPI_ERRORS_RETURN** or a custom MPI Errorhandler

User Level Failure Mitigation

Предложение по смягчению последствий на уровне пользователя (ULFM) было разработано рабочей группой MPI Forum Fault Tolerance [1].

Проектирование механизма, который пользователи будут использовать для управления отказами, строится вокруг трех концепций:

1. **простота**, API должен быть легко понят и использоваться в большинстве распространенных сценариев;
2. **гибкость**, API должен позволять создавать различные отказоустойчивые модели как внешние библиотеки;
3. **отсутствие взаимоблокировки**, никакой вызов MPI (точка-точка или коллективный) не может блокироваться неограниченно после сбоя, он должен либо корректно завершиться, либо сообщить об ошибке.

[1] <http://fault-tolerance.org/>

Overview

- ❑ Code Repository: <https://bitbucket.org/icldistcomp/ulfm>
- ❑ Bulding&Running: <http://fault-tolerance.org/ulfm/ulfm-setup/>
- ❑ ULFM proposal ticket: <https://github.com/mpi-forum/mpi-issues/issues/20>
- ❑ MPI4.0 ???
- ❑ ULFM is prototype implementation in OpenMPI

Сборка&Запуск&Тестирование

- ❑ <http://fault-tolerance.org/ulfm/ulfm-setup/>
- ❑ `hg clone ssh://hg@bitbucket.org/icldistcomp/ulfm`
- ❑ `./configure && make && make install`
- ❑ Сборка и установка происходит также как и OpenMPI

You will need to add one additional header to get access to the new functions, just after the include for `mpi.h`:

For C/C++ this looks like:

```
#include <mpi.h>
```

```
#include <mpi-ext.h>
```

Error-handler sample

```
static void verbose_errhandler(MPI_Comm* comm, int* err, ...) {
    int rank, size, len;
    char errstr[MPI_MAX_ERROR_STRING];
    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
    MPI_Comm_size(MPI_COMM_WORLD, &size);
    MPI_Error_string( *err, errstr, &len );
    printf("Rank %d / %d: Notified of error %s\n", rank, size, errstr);
}

int main(int argc, char *argv[]) {
    int rank, size;
    MPI_Errhandler errh;
    MPI_Init(NULL, NULL);
    MPI_Comm_rank(MPI_COMM_WORLD, &rank); MPI_Comm_size(MPI_COMM_WORLD, &size);
    MPI_Comm_create_errhandler(verbose_errhandler, &errh); // создаем обработчик
    MPI_Comm_set_errhandler(MPI_COMM_WORLD, errh);
    if (rank == (size-1)) { raise(SIGKILL); } // Убиваем процесс
    MPI_Barrier(MPI_COMM_WORLD); // Точка синхронизации
    printf("Rank %d / %d: Stayin' alive!\n", rank, size);
    MPI_Finalize();
    return 0;
}
```

Respawn sample

```
static int MPIX_Comm_replace(MPI_Comm comm, MPI_Comm *newcomm) {
    MPIX_Comm_shrink(comm, &scomm); MPI_Comm_size(scomm, &ns); MPI_Comm_size(comm, &nc);
    MPI_Comm_set_errhandler( scomm, MPI_ERRORS_RETURN );
    rc = MPI_Comm_spawn(...);
}

static int app_needs_repair(MPI_Comm comm) {
    if (comm == world) {
        worldi = (worldi + 1) % 2;
        if (MPI_COMM_NULL != world) { MPI_Comm_free(&world); }
        MPIX_Comm_replace(comm, &world);
        if (world == comm) {
            return false; // Ok, we repaired nothing, no need to redo any work
        }
    }
    return true; // We have repaired the world, we need to reexecute
}

static void errhandler_respawn(MPI_Comm* pcomm, int* errcode, ...) {
    if (MPIX_ERR_PROC_FAILED != eclass && MPIX_ERR_REVOKED != eclass) {
        MPI_Abort(MPI_COMM_WORLD, *errcode);
    }
    MPIX_Comm_revoke(*pcomm);
    app_needs_repair(*pcomm);
}

int main( int argc, char* argv[] ) {
    MPI_Comm_create_errhandler(&errhandler_respawn, &errh);
    while(iteration < max_iterations || app_needs_repair(world)) {
        if (iteration == error_iteration && victim) raise( SIGKILL );
        rc = MPI_Bcast(array, COUNT, MPI_DOUBLE, 0, world);
    }
}
```


Базовые примеры

<http://fault-tolerance.org/ulfm/usage-guide/>

- ❑ **Master/Worker** - The example below presents a master code that handles failures by ignoring failed processes and resubmitting requests. It demonstrates the different failure cases that may occur when posting receptions from `MPI_ANY_SOURCE` as discussed in the advice to users in the proposal.
- ❑ **Iterative Refinement** - The example below demonstrates a method of fault-tolerance to detect and handle failures. At each iteration, the algorithm checks the return code of the `MPI_ALLREDUCE`. If the return code indicates a process failure for at least one process, the algorithm revokes the communicator, agrees on the presence of failures, and later shrinks it to create a new communicator.

```

int master(void) {
    MPI_Comm_set_errhandler(comm, MPI_ERRORS_RETURN);
    MPI_Comm_size(comm, &size);
    // ... submit the initial work requests ...
    MPI_Irecv(buffer, 1, MPI_INT, MPI_ANY_SOURCE, tag, comm, &req);
    // Progress engine: Get answers, send new requests,
    // and handle process failures
    while((active_workers > 0) && work_available) {
        rc = MPI_Wait(&req, &status); // ТОЧКА СИНХРОНИЗАЦИИ!
        if((MPI_ERR_PROC_FAILED == rc) || (MPI_ERR_PENDING == rc)) {
            MPI_Comm_failure_ack(comm); MPI_Comm_failure_get_acked(comm, &g);
            MPI_Group_size(g, &gsize);
            // ... find the lost work and requeue it ...
            active_workers = size - gsize - 1;
            MPI_Group_free(&g);
            // repost the request if it matched the failed process
            if(rc == MPI_ERR_PROC_FAILED) {
                MPI_Irecv(buffer, 1, MPI_INT, MPI_ANY_SOURCE, tag, comm, &req);
            }
            continue;
        }
        // ... process the answer and update work_available ...
        MPI_Irecv(buffer, 1, MPI_INT, MPI_ANY_SOURCE, tag, comm, &req);
    }
    // ... cancel request and cleanup ...
}

```

FAULT TOLERANT MPI APPLICATIONS WITH ULFM

George Bosilca, UTK
Keita Teranishi, SNL
Marc Gamell, Rutgers
Tsutomu Ikegami, AIST
Sara Salem Hamouda, ANU

And many more

SC'15 BoF
Austin, TX, USA



ULFM-based Applications

- **ORNL**: Molecular Dynamic simulation, C/R in memory with Shrink
- **UAB**: transactional FT programming model
- **Tsukuba**: Phalanx Master-worker framework
- **Georgia University**: Wang Landau Polymer Freezing and Collapse, localized subdomain C/R restart
- **Sandia, INRIA, Cray**: PDE sparse solver
- **Cray**: CREST miniapps, PDE solver Schwartz, PPStee (Mesh, automotive), HemeLB (Lattice Boltzmann)
- **ETH Zurich**: Monte-Carlo, on failure the global communicator (that contains spares) is shrunk, ranks reordered to recreate the same domain decomposition
- Programming models: resilient X10

Credits: ETH Zurich

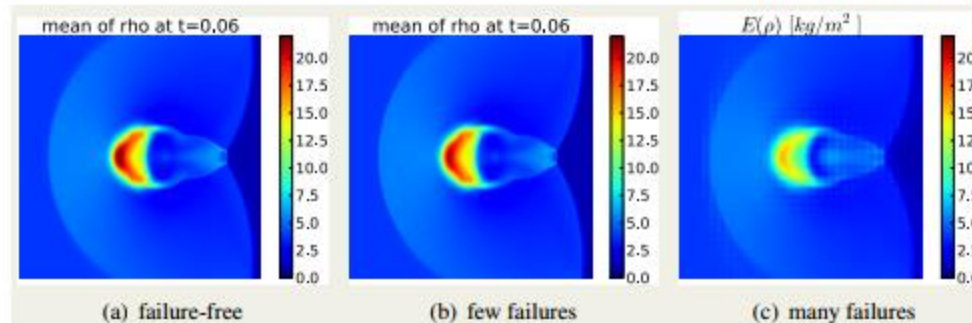
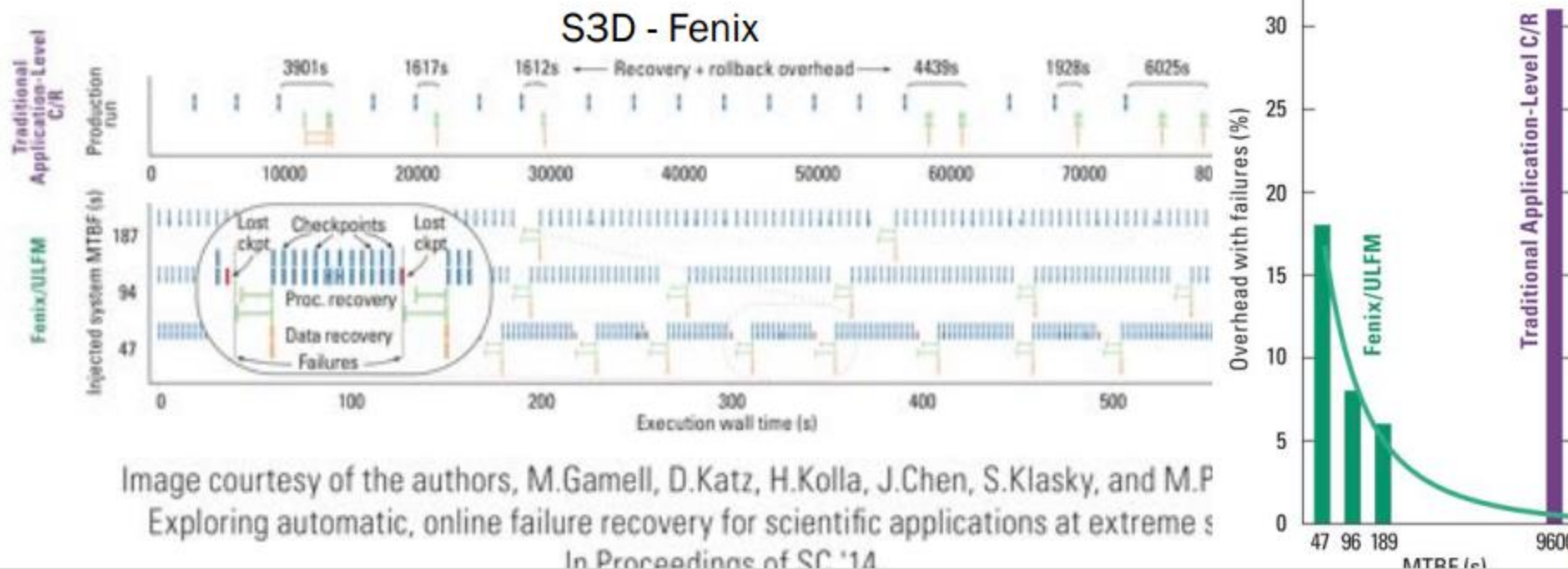


Figure 5. Results of the FT-MLMC implementation for three different failure scenarios.

FRAMEWORKS USING ULFM

LFLR, FENIX, FTLA, Falanx



User activities

- ORNL: Molecular Dynamic simulation
 - Employs coordinated user-level C/R, in place restart with Shrink
- UAB: transactional FT programming model
- Tsukuba: Phalanx Master-worker framework
- Georgia University: Wang Landau Polymer Freezing and Collapse
 - Employs two-level communication scheme with group checkpoints
 - Upon failure, the tightly coupled group restarts from checkpoint, the other distant groups continue undisturbed
- Sandia: PDE sparse solver
- INRIA: Sparse PDE solver
- Cray: CREST miniapps, PDE solver Schwartz, PPStee (Mesh, automotive), HemeLB (Lattice Boltzmann)
- UTK: FTLA (dense Linear Algebra)
 - Employs ABFT
 - FTQR returns an error to the app, App calls new BLACS repair constructs (spawn new processes with MPI_COMM_SPAWN), and re-enters FTQR to resume (ABFT recovery embedded)
- ETH Zurich: Monte-Carlo
 - Upon failure, shrink the global communicator (that contains spares) to recreate the same domain decomposition, restart MC with same rank mapping as before

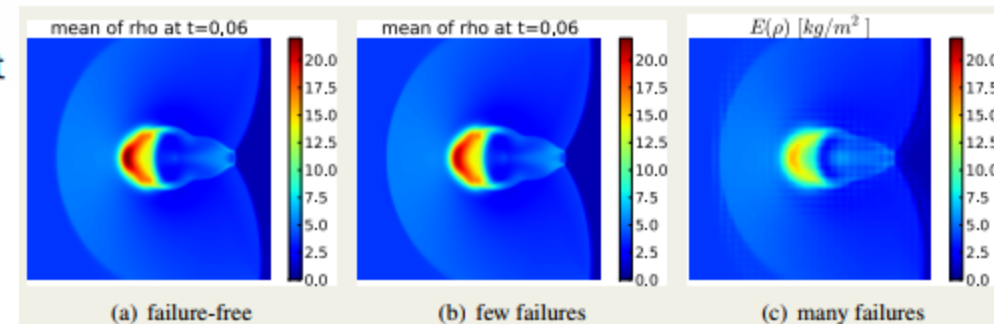
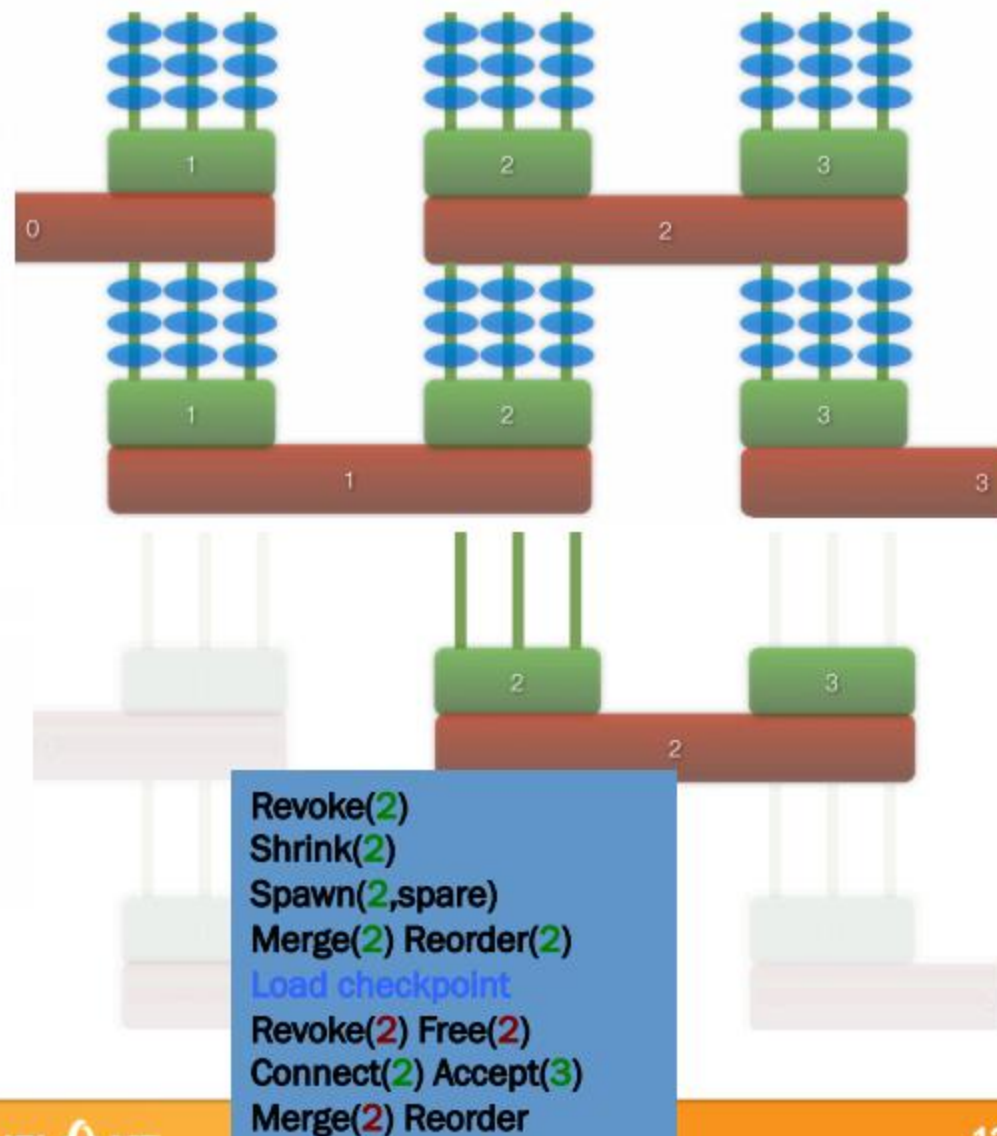


Figure 5. Results of the FT-MLMC implementation for three different failure scenarios.

Credits: ETH Zurich

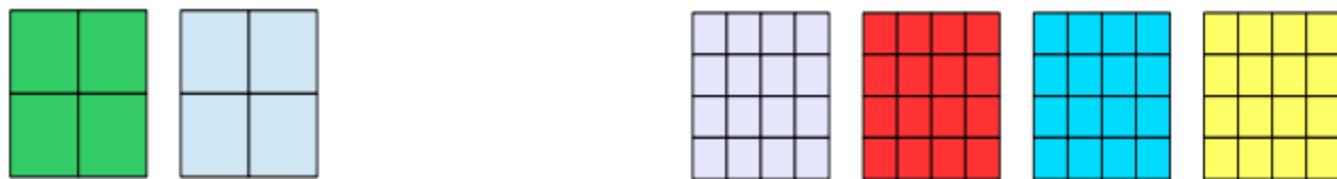
U-GA: Wang-Landau polymer Freeze

- Long independent computation on each processor
 - Dataset protected by **small, cheap checkpoints** (stored on neighbors)
- Periodically, an **AllReduce on the communicator of the Energy window**
- Immediately after, a **Scatter and many pt2pt** on the communicator linking neighboring energy windows



ETH-Zurich: Monte-Carlo PDE

- X is the solution to a stochastic PDE
- Each sample X^i is computed with a FVM solver
- MC error is determined by
 - stochastic error (depends on M)
 - discretization error (depends on the mesh-width h)
- A more accurate MC approximation requires more samples M and a finer mesh h

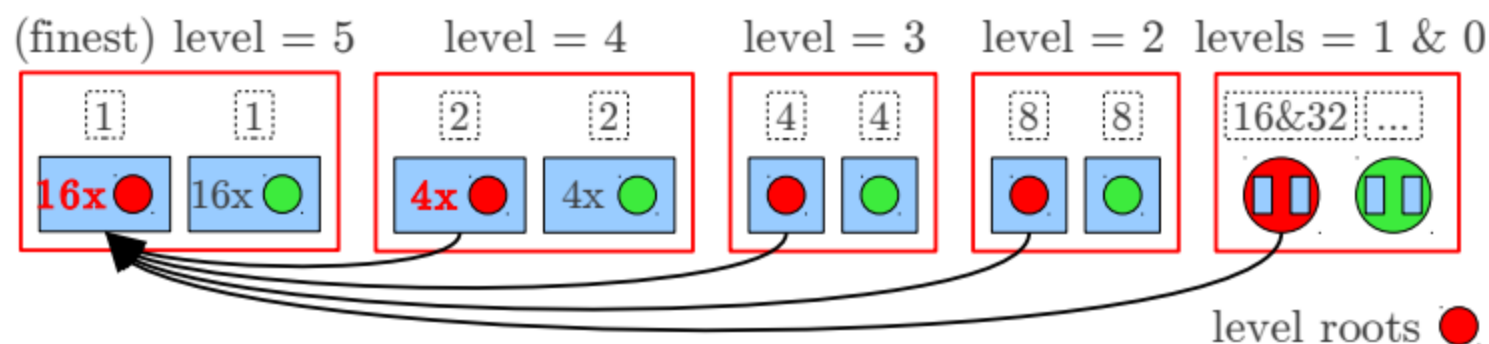


Fault Tolerant Monte Carlo:

- The number of samples M turns into a random variable \hat{M}

- $$\|\mathbb{E}[X] - E_{\hat{M}}[X]\| \leq \mathbb{E} \left[\frac{1}{\sqrt{\hat{M}}} \right] \|X\|$$

ETH-Zurich: Monte-Carlo PDE



- Try to collect the mean as in fault-free ALSVID-UQ
- Call `MPI_BARRIER` on `MPI_COMM_WORLD` at the end to discover failed processes
- non-uniform success of `MPI_BARRIER`: `MPI_BARRIER` is followed by `MPI_COMM_AGREE`
- In case of failure: (Re)assign the level roots and repeat the collection of the means

Simplifying the Recovery Model of User-Level Failure Mitigation

Wesley Bland

ExaMPI '14

New Orleans, LA, USA

November 17, 2014

What has been done before?

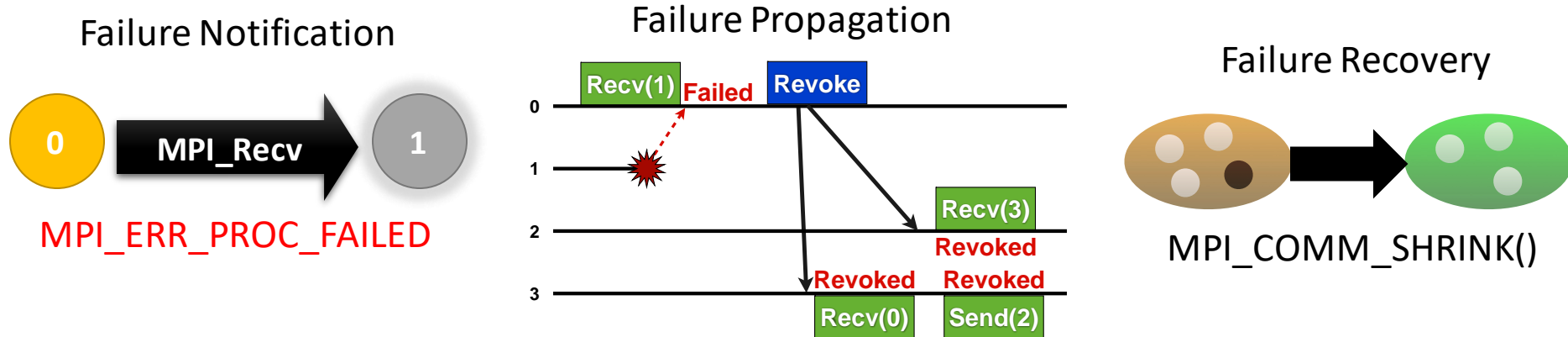
- Checkpoint/Restart
 - Requires restarting the entire job, waiting in job queue
- FT-MPI (University of Tennessee)
 - Automatic communicator repair, no customization
 - Support dropped and never presented to the MPI Forum
- FA-MPI (University of Alabama, Birmingham)
 - Try/catch semantics to simulate transactions
 - Detects errors via timeouts
- Run-Through Stabilization (MPI Forum)
 - Similar to ULFM, but too big to not fail

User-Level Failure Mitigation (ULFM)

- Handles (fail-stop) process failures
- General resilience framework designed to support a wide range of recovery models
- Encourages libraries to build more user-friendly resilience on top of ULFM
- Failure notification via return codes and/or MPI_Errhandlers
 - Local failure notification to maintain performance

ULFM Continued

- Failure propagation happens manually when necessary
 - **MPI_COMM_REVOKE**
 - All communication operations on a communicator are interrupted
- New function call creates a replacement communicator without failed procs
 - **MPI_COMM_SHRINK**
- Group of failed procs available via API calls



FT Agreement

- Collective over communicator
- Allows procs to determine the status of an entire communicator in one call
 - Bitwise AND over an integer value
- Ignores (acknowledged) failed processes
 - Provides notification for unacknowledged failures
- Works on a revoked communicator
- Useful for “validating” previous work
 - Communicator creation
 - Application iteration(s)

ULFM for BSP

- Take advantage of existing synchrony
- Perform agreement on the result of the ending collective
- If there's a problem, repair the communicator and data and re-execute
- Avoid restarting the job
 - Sitting in the queue, restarting the job, re-reading from disk, etc.

```
while (<still working>) {
```

```
...  
Application Code  
...
```

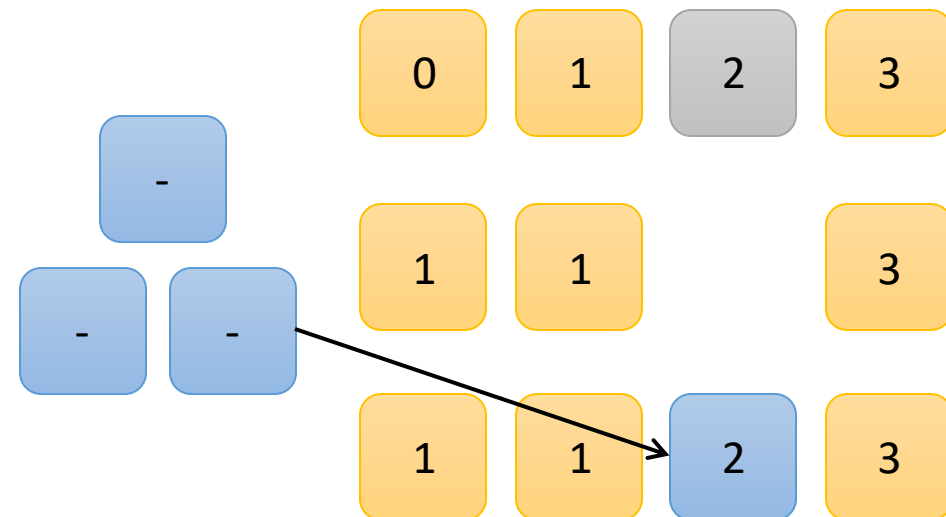
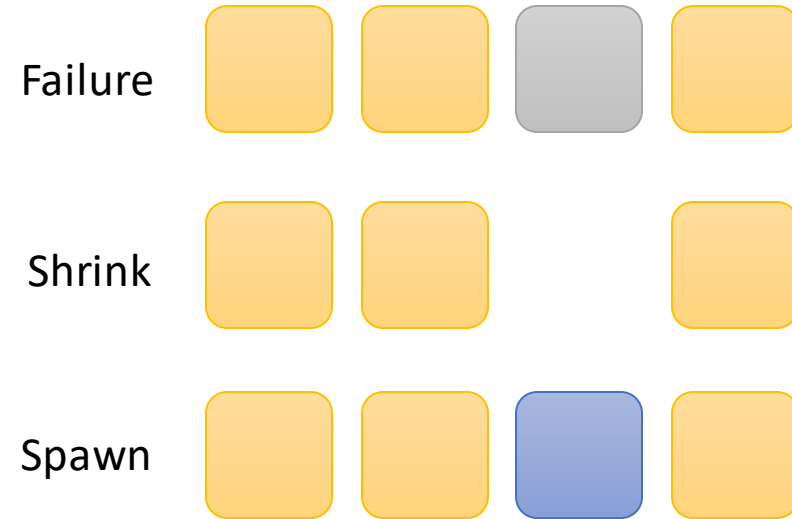
```
rc = MPI_Allreduce(value)
```

```
if (rc == MPI_SUCCESS) success = 1  
rc = MPI_Comm_agree(success)  
if (!success || !rc) repair()
```

```
}
```

Recovery Model

- Communication Recovery
 - Dynamic processes (+ batch support)
 - Shrink + Spawn
 - No-dynamic processes
 - Start the job with extra processes
 - Rearrange after a failure
- Data Recovery
 - Full checkpoint restart
 - Already built into many apps
 - High recovery and checkpointing overhead
 - Application level checkpointing
 - Low overhead
 - Possible code modifications required



Code Intrusiveness

- Add 2-3 lines to check for failures
 - Perform agreement to determine status
- Add function to repair application
 - If the agreement says there's a failure, call the repair function
 - Repair data and communicator(s)

Monte Carlo Communication Kernel (MCCK)

- Mini-app from the Center for Exascale Simulation of Advanced Reactors (CESAR)
- Monte-Carlo, domain decomposition stencil code
- Investigates communication costs of particle tracking

ULFM Modifications for MCCK

Application Level Checkpointing

- Checkpoint particle position to disk
- On failure, restart job
 - No job queue to introduce delay
- All processes restore data from disk together

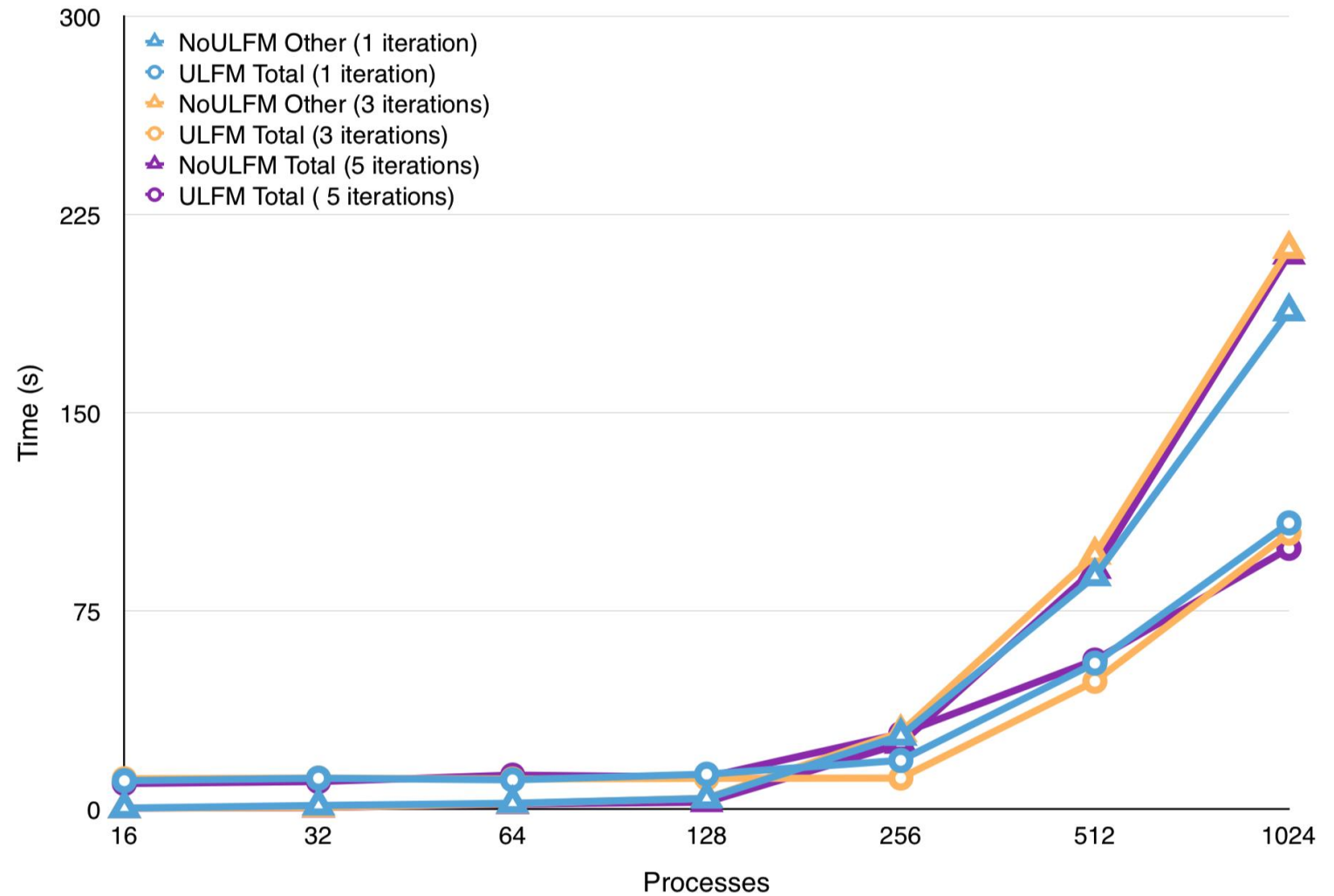
User Level Failure Mitigation

- Particle position exchanged with neighbors
- On failure, call repair function
 - Shrink communicator
 - Spawn replacement process
 - Restore missing data from neighbor
- No disk contention

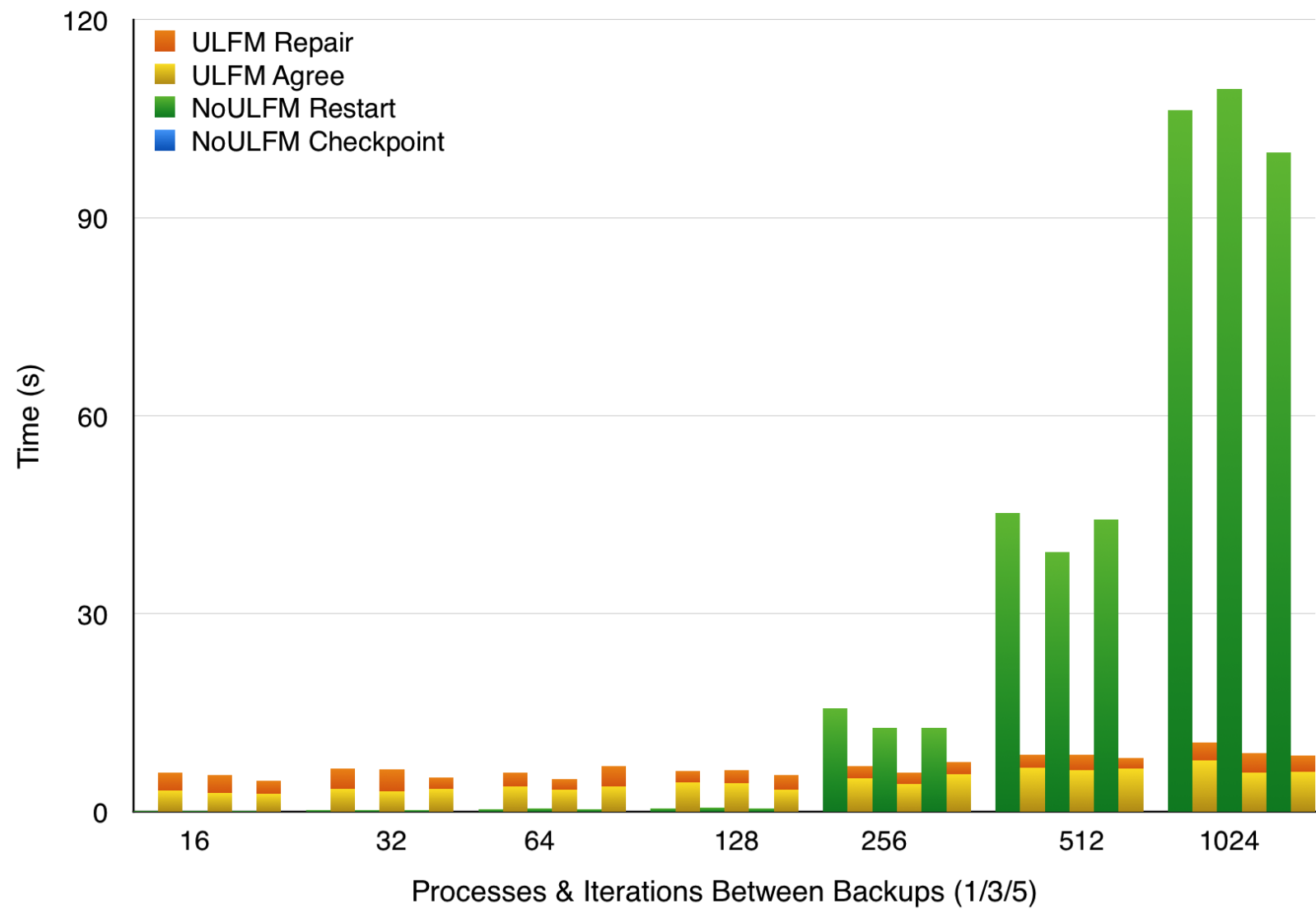
Performance

- Used Fusion cluster @ Argonne
 - 2 Xeons per node (8 cores)
 - 36 GB memory
 - QDR InfiniBand
- Experiments up to 1024 processes
- 100,000 particles per process
- Checkpoints taken every 1, 3, and 5 iterations

Overall Runtime



Overhead



What's next for ULFM?

- Fault tolerance libraries
 - Incorporate data & communication recovery
 - Abstract details from user applications
- Working on standardization in MPI Forum

ULFM описание нового функционала

В данной библиотеке каждому MPI-процессу поставлен в соответствие атрибут, значение которого отражает состояние этого процесса: отказал он или нет.

Также в ней реализованы модификации MPI функций обмена. Функции обмена данной библиотеки возвращают исключения согласно спецификации расширения ULFM, не предусмотренные стандартом MPI 3.0, позволяющие моделировать отказ в нормально функционирующей распределенной вычислительной системе.

Для реализации техник отказоустойчивости с помощью функций стандарта MPI 3.0 в этой библиотеке реализованы функции

**MPI_COMM_REVOKE, MPI_COMM_SHRINK, MPI_COMM_AGREE,
MPI_COMM_FAILURE_ACK, MPI_COMM_FAILURE_GET_ACKED**

***МОДЕЛИРОВАНИЕ ОТКАЗОВ В ВЫСОКОПРОИЗВОДИТЕЛЬНЫХ ВЫЧИСЛИТЕЛЬНЫХ СИСТЕМАХ В РАМКАХ
СТАНДАРТА MPI И ЕГО РАСШИРЕНИЯ ULFM1 А.А. Бондаренко, М.В. Якобовский*

ULFM описание нового функционала

Обработка исключений и реализация различных техник восстановления осуществляется, в основном за счет функций:

- ❑ **MPI_COMM_REVOKE** — прекращает все текущие нелокальные операции на коммуникаторе и отмечает коммуникатор в качестве аннулированного. Все последующие вызовы нелокальных функций, связанные с этим коммуникатором, должны завершаться значением **MPI_ERR_REVOKED**, за исключением функций **MPI_COMM_SHRINK** и **MPI_COMM_AGREE**;
- ❑ **MPI_COMM_SHRINK** — создает на основе существующего коммуникатора новый, не содержащий отказавшие процессы;
- ❑ **MPI_COMM_FAILURE_GET_ACKED** — возвращает группу, состоящую из процессов, которые были определены как отказавшие к моменту последнего вызова функции **MPI_COMM_FAILURE_ACK**;
- ❑ **MPI_COMM_AGREE** — согласовывает значение булевой переменной, если нет отказавших MPI-процессов в коммуникаторе или возвращает исключение о наличии отказа всем не отказавшим процессам в коммуникаторе.

Requirements for MPI standardization of FT

- **Expressive**, simple to use
 - Support legacy code, **backward compatible**
 - Enable users to port their code simply
 - **Support a variety of FT models and approaches**
- Minimal (ideally **zero**) impact on failure free **performance**
 - No global knowledge of failures
 - No supplementary communications to maintain global state
 - Realistic memory requirements
- **Simple to implement**
 - Minimal (or **zero**) changes to existing functions
 - Limited number of new functions
 - Consider thread safety when designing the API

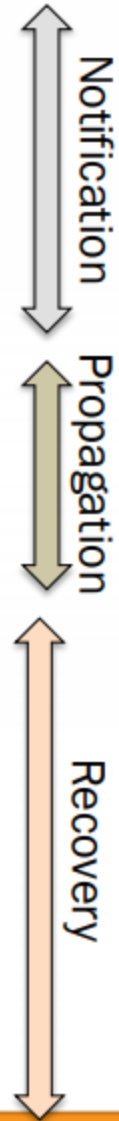


Integration with existing mechanisms

- New error codes to deal with failures
 - **MPI_ERROR_PROC_FAILED**: report that the operation discovered a newly dead process. Returned from all blocking function, and all completion functions.
 - **MPI_ERROR_PROC_FAILED_PENDING**: report that a non-blocking MPI_ANY_SOURCE potential sender has been discovered dead.
 - **MPI_ERROR_REVOKED**: a communicator has been declared improper for further communications. All future communications on this communicator will raise the same error code, with the exception of a handful of recovery functions
- Is that all?
 - Matching order (MPI_ANY_SOURCE), collective communications

Summary of new functions

- **MPI_Comm_failure_ack(comm)**
 - Resumes matching for MPI_ANY_SOURCE
- **MPI_Comm_failure_get_acked(comm, &group)**
 - Returns to the user the group of processes acknowledged to have failed
- **MPI_Comm_revoke(comm)**
 - **Non-collective** collective, interrupts all operations on comm (future or active, at all ranks) by raising MPI_ERR_REVOKED
- **MPI_Comm_shrink(comm, &newcomm)**
 - Collective, creates a new communicator without failed processes (identical at all ranks)
- **MPI_Comm_agree(comm, &mask)**
 - Collective, agrees on the AND value on binary mask, ignoring failed processes (reliable AllReduce), and the return core



MPI_Comm_failure_ack

- **Local** operations that acknowledge all locally notified failures
 - Updates the group returned by MPI_COMM_FAILURE_GET_ACKED
- Unmatched MPI_ANY_SOURCE that would have raised MPI_ERR_PROC_FAILED_PENDING proceed without further exceptions regarding the acknowledged failures.
- MPI_COMM_AGREE do not raise MPI_ERR_PROC_FAILED due to acknowledged failures
 - No impact on other MPI calls especially not on collective communications

MPI_Comm_failure_get_acked

- Local operation returning the group of failed processes in the associated communicator that have been locally acknowledged
- Hint: All calls to `MPI_Comm_failure_get_acked` between a set of `MPI_Comm_failure_ack` return the same set of failed processes

Failure Discovery

- Discovery of failures is *local* (different processes may know of different failures)
- **MPI_COMM_FAILURE_ACK(comm)**
 - This local operation gives the users a way to acknowledge all locally notified failures on comm. After the call, unmatched MPI_ANY_SOURCE receive operations proceed without further raising MPI_ERR_PROC_FAILED_PENDING due to those acknowledged failures.
- **MPI_COMM_FAILURE_GET_ACKED(comm, &grp)**
 - This local operation returns the group *grp* of processes, from the communicator comm, that have been locally acknowledged as failed by preceding calls to MPI_COMM_FAILURE_ACK.
- Employing the combination ack/get_acked, a process can obtain the list of all failed ranks (as seen from its local perspective)

MPI_Comm_revoke

- Communicator level failure propagation
- The revocation of a communicator complete all pending local operations
 - A communicator is revoked either after a local MPI_Comm_revoke or any MPI call raise an exception of class MPI_ERR_REVOKED
- Unlike any other concept in MPI it is not a collective call but has a collective scope
- Once a communicator has been revoked all non-local calls are considered local and must complete by raising MPI_ERR_REVOKED
 - Notable exceptions: the recovery functions (agreement and shrink)

MPI_Comm_revoke

- Communicator level failure propagation
- The revocation of a communicator complete all pending local operations
 - A communicator is revoked either after a local MPI_Comm_revoke or any MPI call raise an exception of class MPI_ERR_REVOKED
- Unlike any other concept in MPI it is not a collective call but has a collective scope
- Once a communicator has been revoked all non-local calls are considered local and must complete by raising MPI_ERR_REVOKED
 - Notable exceptions: the recovery functions (agreement and shrink)

MPI_Comm_agree

- Perform a consensus between all living processes in the associated communicator and consistently return a value and an error code to all living processes
- Upon completion all living processes agree to set the output integer value to a bitwise AND operation over all the contributed values
 - Also perform a consensus on the set of known failed processes (!)
 - Failures non acknowledged by all participants raise MPI_ERR_PROC_FAILED

MPI_Comm_shrink

- Creates a new communicator by excluding all known failed processes from the parent communicator
 - It completes an agreement on the parent communicator
 - Work on revoked communicators as a mean to create safe, globally consistent sub-communicators
- Absorbs new failures, it is not allowed to return MPI_ERR_PROC_FAILED or MPI_ERR_REVOKED