# A synthetic tool for analysing adaptive workloads

*Authors*:
<u>Iker Martín-Álvarez</u>
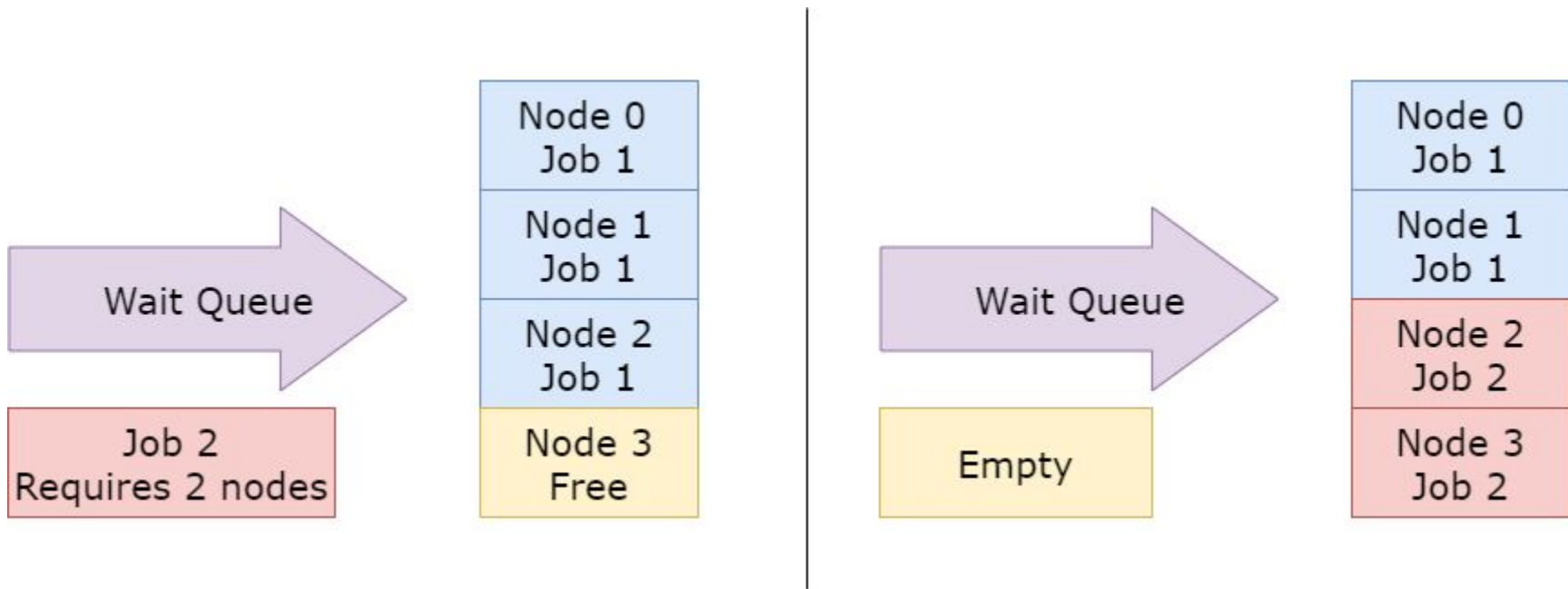José I. Aliaga
María Isabel Castillo
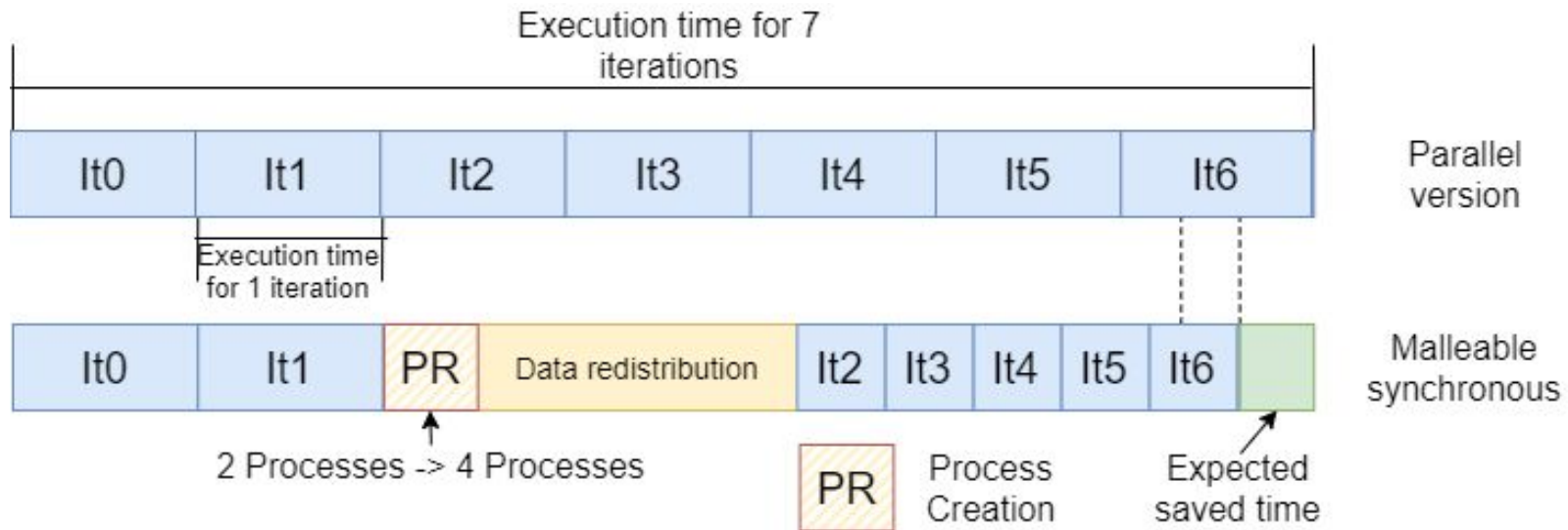Sergio Iserte
Rafael Mayo

Outline:

- Motivation
- Synthetic application
    - Computation module
    - Reconfiguration module
    - Methric gatherer module
    - Configuration file parameters
- Results
    - Total execution times
    - Malleability times
    - Iteration times
- Conclusions
- Future work

Execution time for 7 iterations

| It0 | It1 | It2 | It3 | It4 | It5 | It6 | Parallel version |

Execution time for 1 iteration

| It0 | It1 | PR | Data redistribution | It2 | It3 | It4 | It5 | It6 | | Malleable synchronous |

2 Processes -> 4 Processes

PR — Process Creation

Expected saved time

Motivation (II)

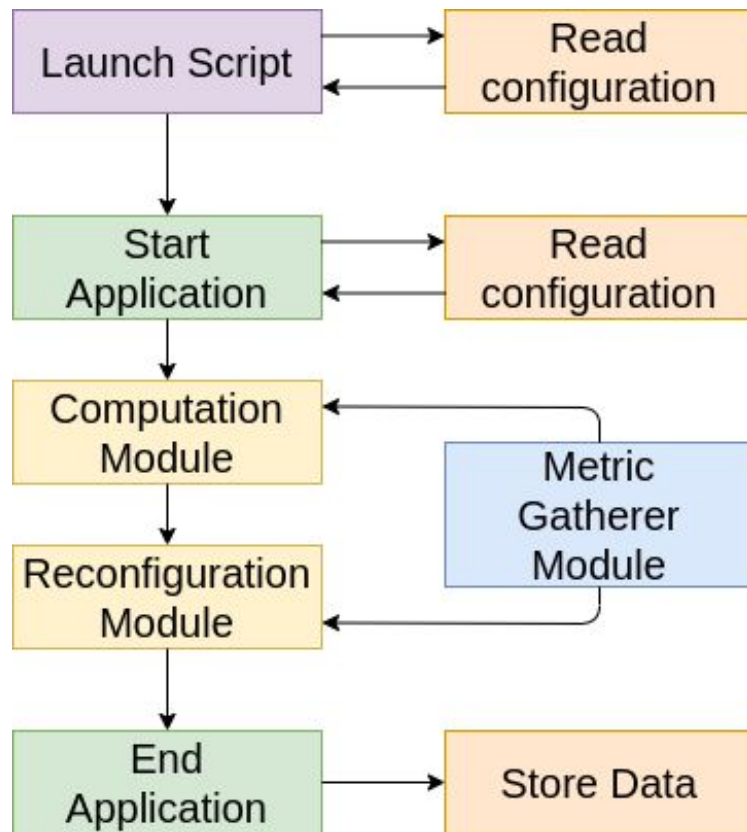UNIVERSITAT JAUME I

# Synthetic application

Application composed of three modules:
- Computation module
- Reconfiguration module
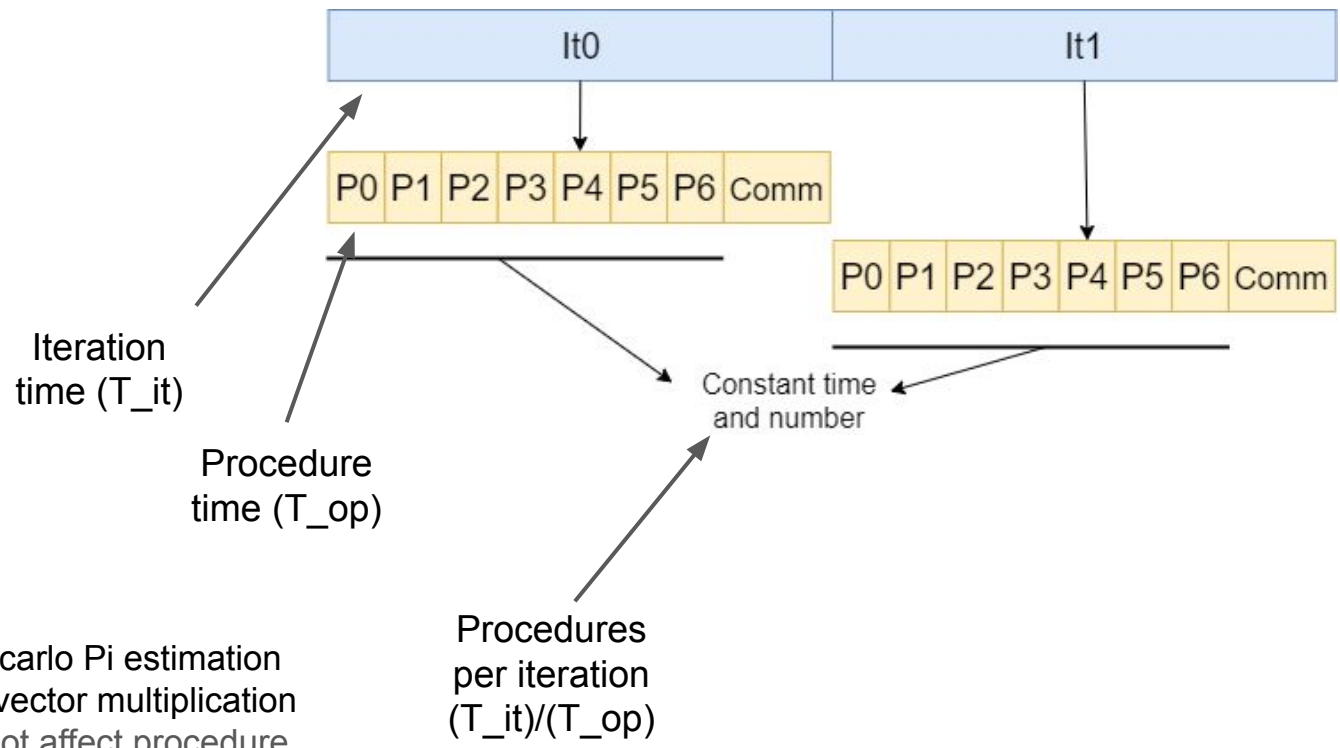- Metric gatherer module

and a configuration file

```
int main() {
  // First group initializes the application
  if(init) {
    if (rank == 0)
      read_config();
    broadcast_config(); // Send to all ranks in group
    calculate_non_direct_parameters();
  } else {
    broadcast_config() // Receive from parent rank 0
    calculate_non_direct_parameters();
    redistribute_data(); //receive from parents.
  }
  // Job computation
  for(int iter=0; iter < Iters; iter++) {
    // Computes for consume T_it time
    for(int i=0; i < op; i++) {
      compute(N, Pt);
    }
    // Communicates Cb bytes (optional)
    communications(Cb, Ct);
  }
  // Reconfigure if not last resize
  if (!last_resize) {
    create_child_processes();
    broadcast_config();  // From rank 0 to all children
    redistribute_data(); //send to children
  }
  store_performance_data();
}
```

# Computation module
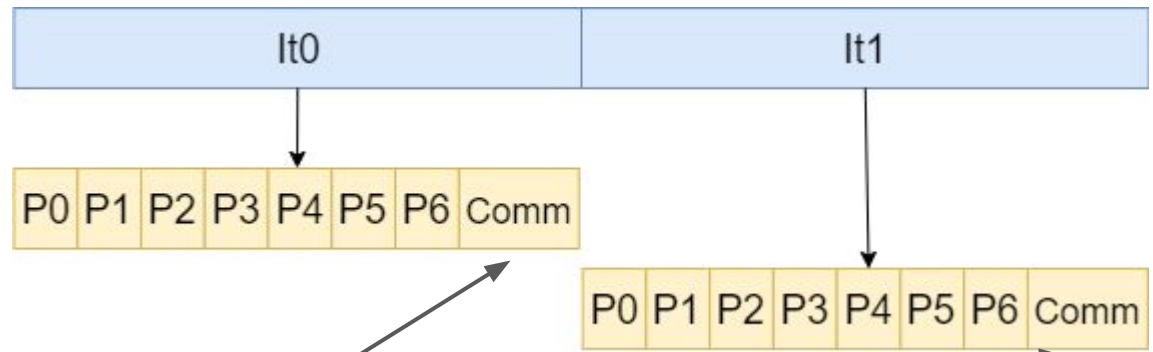
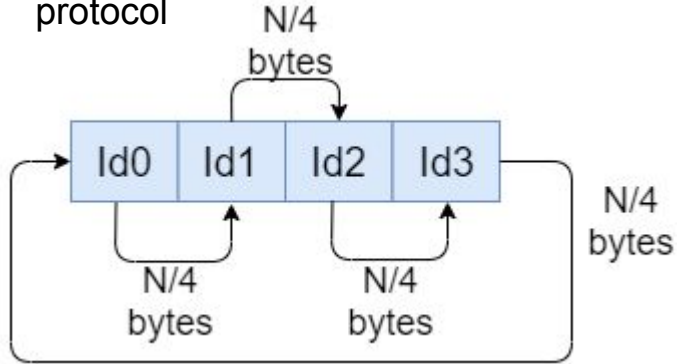Iteration time (T_it)

Procedure time (T_op)

Computation procedures:
- Compute-bound: Montecarlo Pi estimation
- Memory-bound: Matrix-vector multiplication
- Process number does not affect procedure final time (T_op)

Procedures per iteration (T_it)/(T_op)

Communication procedures:
- Point to point (MPI_Send/Recv)
- Collective one-to-all (MPI_Bcast)
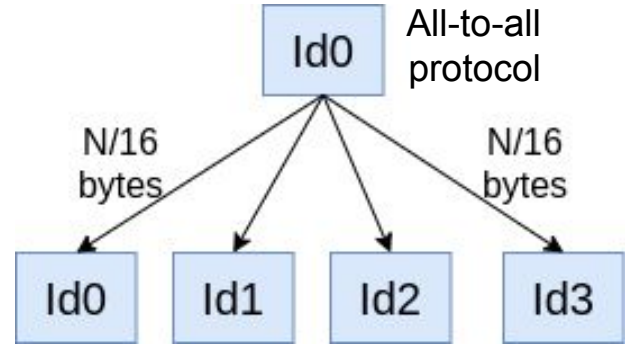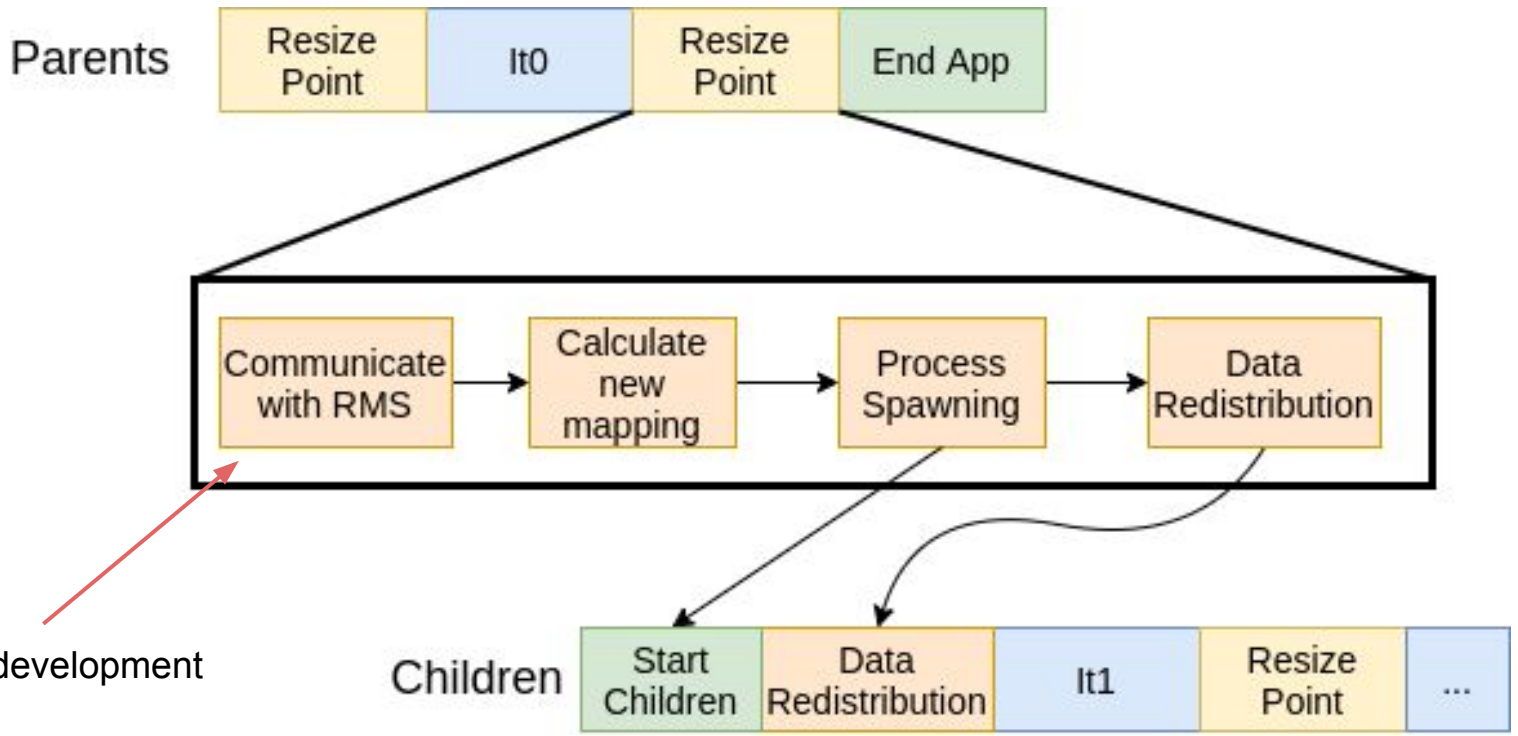- Collective all-to-all (MPI_Alltoall)
- Reduction (MPI_Reduce)

| It0 | It1 |
|-----|-----|

| P0 | P1 | P2 | P3 | P4 | P5 | P6 | Comm |
|----|----|----|----|----|----|----|------|

| P0 | P1 | P2 | P3 | P4 | P5 | P6 | Comm |
|----|----|----|----|----|----|----|------|

Communications cost:
- *N* bytes

Point to point protocol

N/4 bytes

| Id0 | Id1 | Id2 | Id3 |
|-----|-----|-----|-----|

N/4 bytes

N/4 bytes

N/4 bytes

N/4 bytes

All-to-all protocol

Id0

N/16 bytes

N/16 bytes

| Id0 | Id1 | Id2 | Id3 |
|-----|-----|-----|-----|

# Reconfiguration module

Reconfiguration module - Steps (I)

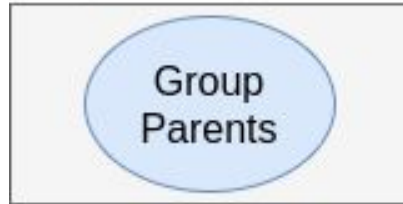job_info:
- Node number
- Core number
- Node names

Under development

Forced resize

Parents: It0 | It1 | Resize Point | End App

Calculate new mapping → Process Spawning → Data Redistribution

Children: Start Children | Data Redistribution | It2 | ...

Reconfiguration module - Physical mappings

Reconfiguration module - Process spawn

Parents

| It0 | It1 | PR | Data redistribution | |

Data redistribution | It2 | It3 | It4 | It5 | It6

Children

Based on MPI_Alltoallv

Distribution of elements from parent P0 to its children C0, C1 and C2

Distribution from 4 to 10 processes of 100 elements

UJI UNIVERSITAT JAUME I   Data redistribution - Basics

**Additional contribution: Allow asynchronous data redistribution**

SR = Synchronous redistribution

AR = Asynchronous redistribution
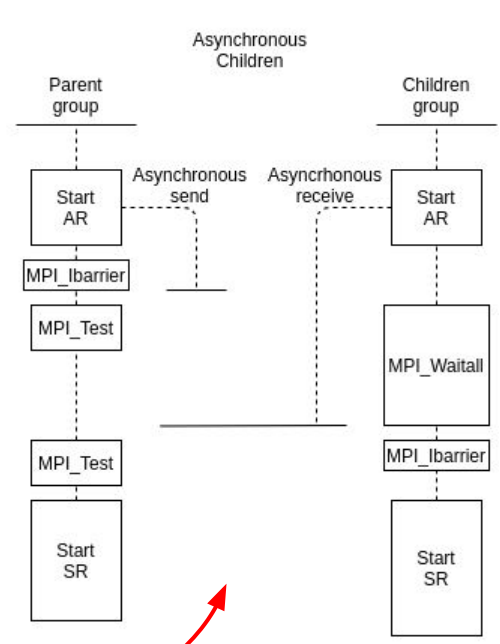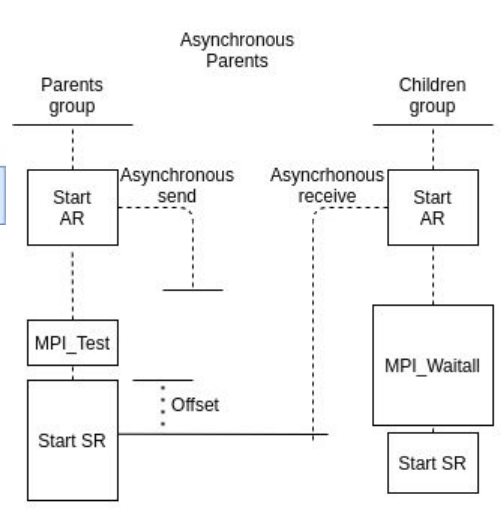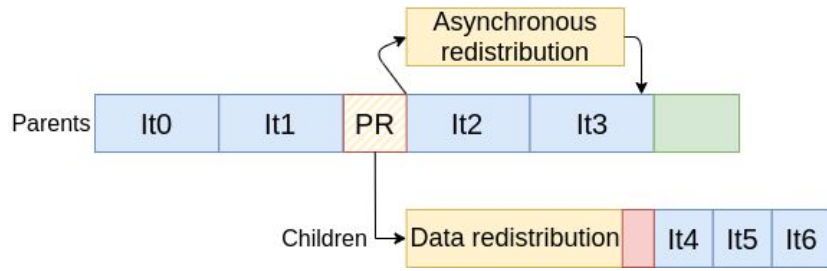
Execution time for 7 iterations

| It0 | It1 | It2 | It3 | It4 | It5 | It6 | Parallel version |

Execution time for 1 iteration

| It0 | It1 | PR | Data redistribution | It2 | It3 | It4 | It5 | It6 | Malleable synchronous |

2 Processes -> 4 Processes

| It0 | It1 | PR | It2 | It3 | It4 | It5 | It6 | Malleable asynchronous |

PR — Process Creation

Asynchronous redistribution

Expected saved time

UNIVERSITAT JAUME I

Data redistribution - Asynchronous (I)

Asynchronous procedures:
- Point to point (MPI_Isend/Irecv)
- Collective all-to-all Parents (MPI_Ialltoallv)
- Collective all-to-all Children (MPI_Ialltoallv)
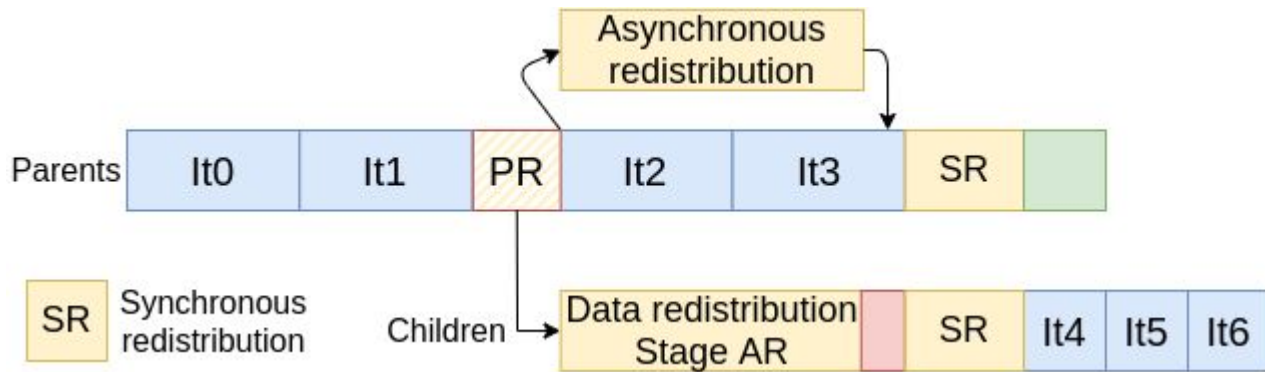- Pthreads (MPI_Alltoallv)

MPI Non-blocking Primitives

Asynchronous procedures:
- Point to point (MPI_Isend/Irecv)
- Collective all-to-all Parents (MPI_Ialltoallv)
- Collective all-to-all Children (MPI_Ialltoallv)
- Pthreads (MPI_Alltoallv)

Parents: It0 | It1 | PR | It2 | It3 | SR

Asynchronous redistribution

SR — Synchronous redistribution

Children → Data redistribution Stage AR | SR | It4 | It5 | It6

Stage 1

Stage 2

| Type | Communication mechanism |
|---|---|
| Constant | Synchronous / Asynchronous |
| Variable | Strictly Synchronous |

# Metric gatherer module

Based on MPI_Wtime

Total execution time

| It0 | It1 | PR | Data redistribution | It2 | It3 | It4 | It5 | It6 |

Synchronous redistribution time

Process spawn time

Iteration time

Also stored with each iteration:
- Number of *Op* executed
- If AR is being performed or not

| It0 | It1 | PR | It2 | It3 | It4 | It5 | It6 |

Asynchronous redistribution

Asynchronous redistribution time

UNIVERSITAT JAUME I

Metric gatherer module

# Configuration files

Configuration file - General

```
R=1                  # Number of resizes
SDR=1000000000       # Synchronously redistributed bytes
ADR=1000000000       # Asynchronously redistributed bytes
AT=3                 # Type of asynchronous redistribution
;end [general]
```

```
[resize0]
Iters=1          # Total iterations for this group
Procs=2          # Number of processes in this group
FactorS=0.5      # Scalability factor
Dist=balanced    # Physical process mapping type
;end [resize0]
[resize1]
Iters=10
Procs=10
FactorS=0.1
Dist=balanced
;end [resize1]
```

| It0 | PR | Data redistribution | It1 | ... | It10 |

```
[resize0]
Iters=1          # Total iterations for this group
Procs=2          # Number of processes in this group
FactorS=0.5      # Scalability factor
Dist=balanced    # Physical process mapping type
;end [resize0]
[resize1]
Iters=10
Procs=10
FactorS=0.1
Dist=balanced
;end [resize1]
```

$Real\_It = T\_it * FactorS$

$Real\_It = 4 * 0.5 = 2s$

| It0 | PR | Data redistribution | It1 | ... | It10 |

$Real\_It = 4 * 0.1 = 0.4s$

Application composed of three modules:
- Computation module
- Reconfiguration module
- Metric gatherer module

and a configuration file

```
int main() {
  // First group initializes the application
  if(init) {
    if (rank == 0)
      read_config();
    broadcast_config(); // Send to all ranks in group
    calculate_non_direct_parameters();
  } else {
    broadcast_config() // Receive from parent rank 0
    calculate_non_direct_parameters();
    redistribute_data(); //receive from parents.
  }
  // Job computation
  for(int iter=0; iter < Iters; iter++) {
    // Computes for consume T_it time
    for(int i=0; i < op; i++) {
      compute(N, Pt);
    }
    // Communicates Cb bytes (optional)
    communications(Cb, Ct);
  }
  // Reconfigure if not last resize
  if (!last_resize) {
    create_child_processes();
    broadcast_config();  // From rank 0 to all children
    redistribute_data(); //send to children
  }
  store_performance_data();
}
```

# Results

- Two 10-core processor Intel Xeon 4210
- Two servers (40 cores)
- 11 Iterations of 4s or 0.4s per execution
- 10MB iteration communication
- One reconfiguration per execution
- *FactorS* is applied to have a perfect SpeedUp
- Both physical mappings
- Fully SR or AR of 1GB
- Pthreads option
- Mean of 5 executions

Compact mapping

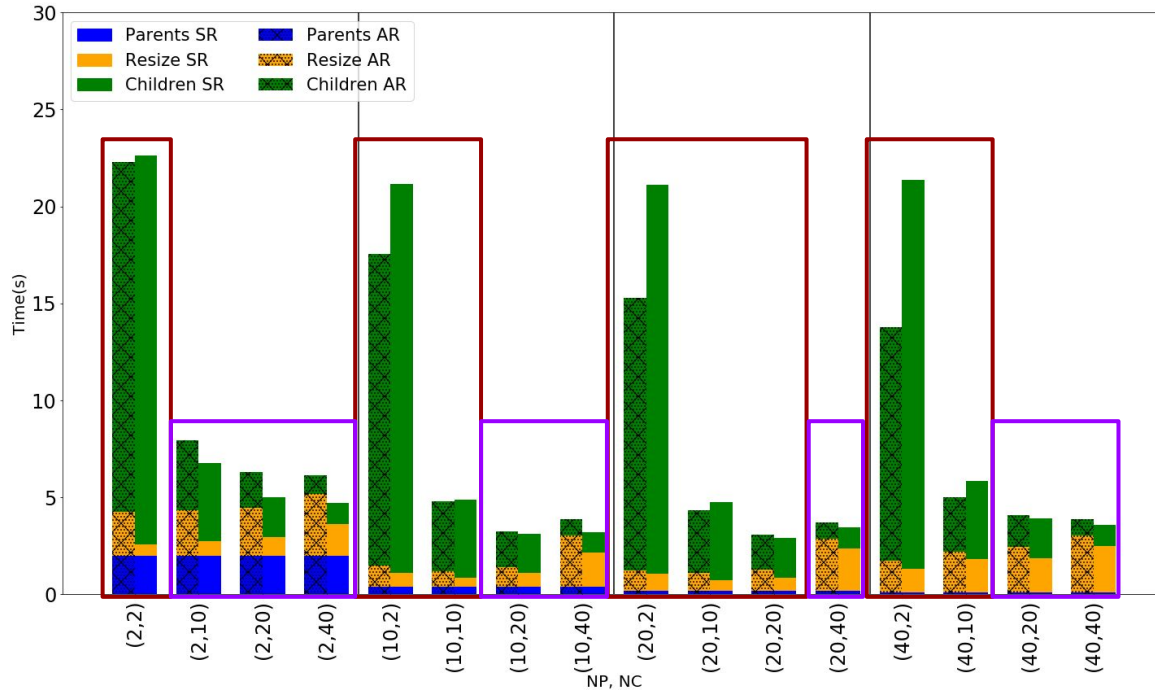T_it = 4 seconds

Shrinking with AR

Expanding with SR

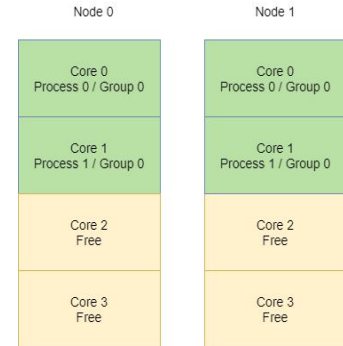80 processes, 120 threads in AR

Oversubscription when more than 20 threads in one node.

Results - Total execution times (I)

Balanced mapping

T_it = 4 seconds

Shrinking with AR

Expanding with SR

Oversubscription when more than 20 threads.

Compact mapping

T_it = 0.4 seconds

Shrinking with AR

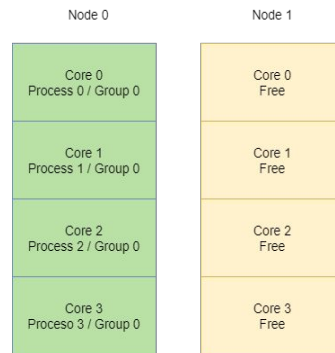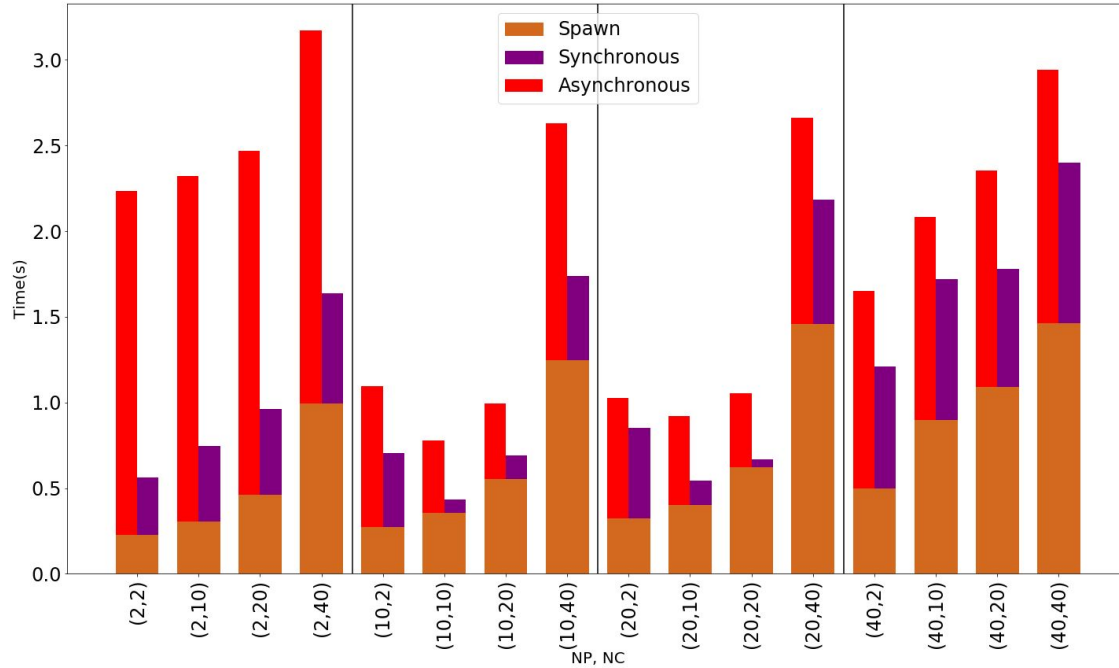Expanding with SR

Oversubscription when more than 20 threads.

Results - Final execution times (III)

UNIVERSITAT JAUME I

Balanced mapping          T_it = 0.4 seconds

Shrinking with AR

Expanding with SR

Oversubscription when more than 20 threads.

Results - Final execution times (IV)

Compact mapping          T_it = 4 seconds
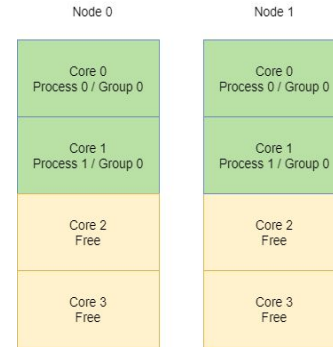
Spawning processes is the most expensive operation

AR is always more expensive than SR

Results - Reconfiguration times (I)

Balanced mapping        T_it = 4 seconds
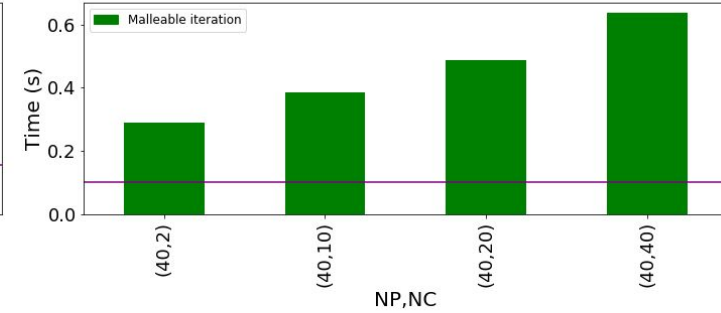
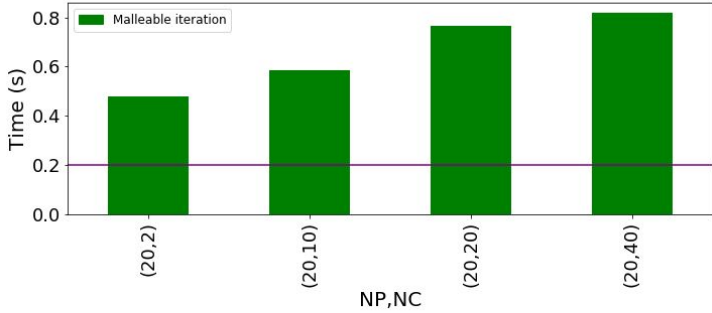Spawning processes is the most expensive operation

AR is always more expensive than SR

Better performance than Cm: Oversubscription appears in less configurations

Results - Reconfiguration times (II)

UNIVERSITAT JAUME I

Compact mapping    T_it = 0.4 seconds

Spawning processes is the most expensive operation

AR is always more expensive than SR

AR is better for lower T_it values

Balanced mapping          T_it = 0.4 seconds

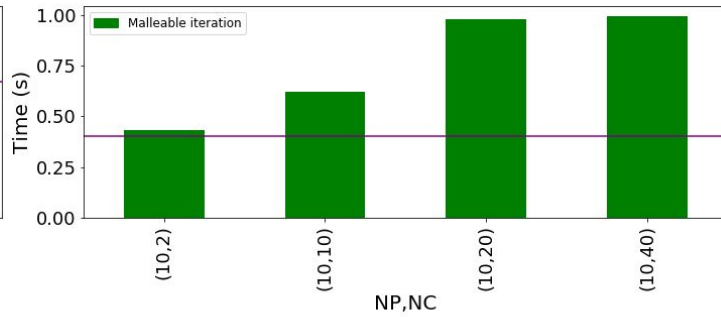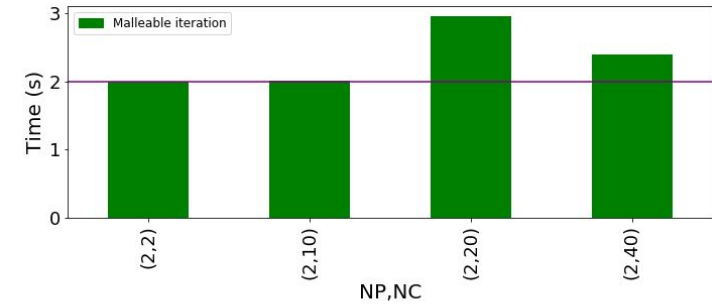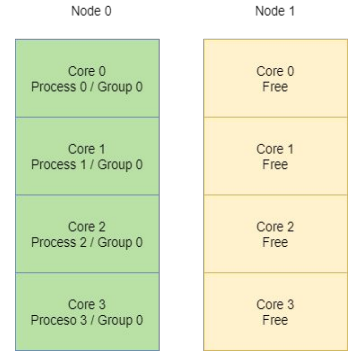Spawning processes is the most expensive operation

AR is always more expensive than SR

Better performance than Cm: Oversubscription appears in less configurations
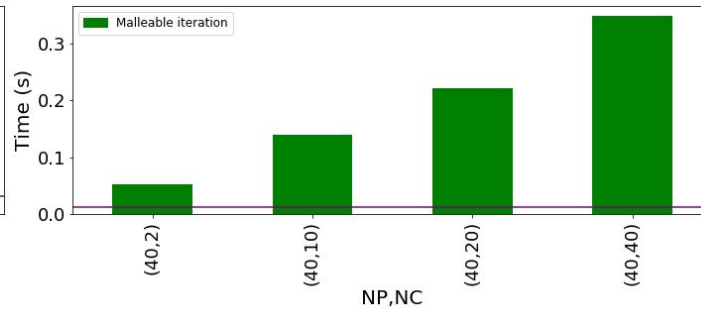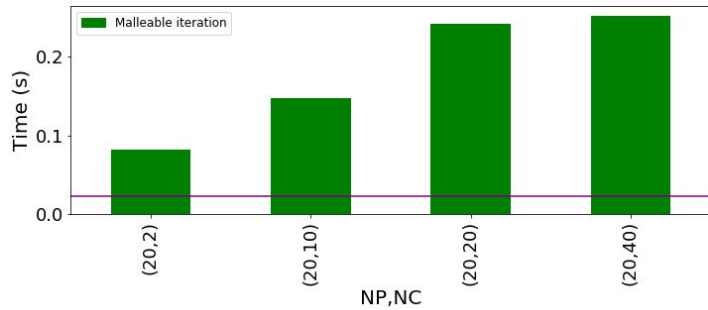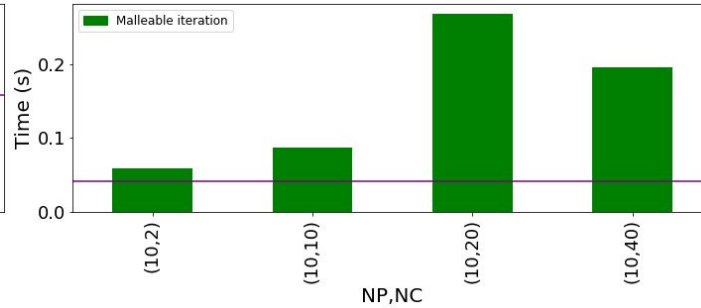
AR is better for lower T_it values

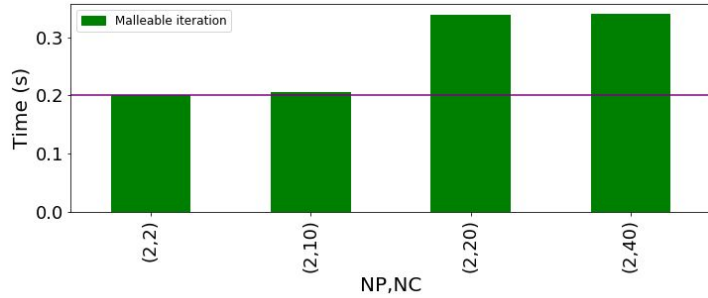Results - Reconfiguration times (IV)

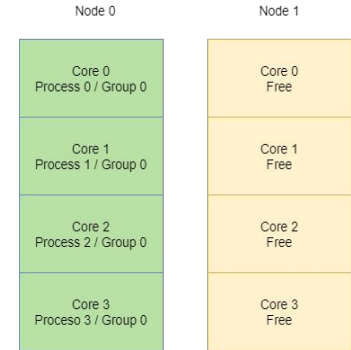UNIVERSITAT JAUME I

Compact mapping     T_it = 4 seconds

Oversubscription reduces iteration performance drastically

Compact mapping       T_it = 0.4 seconds

Oversubscription reduces iteration performance drastically

Results - Iteration times under reconfiguration (II)

UNIVERSITAT JAUME I

**Adaptable synthetic application for expanding, shrinking or migrating**

**Application allows to study which reconfiguration mechanism is preferred depending on the job state**

**SR when expanding the job, AR when shrinking**

**Oversubscription reduces performance for AR and SR**

UNIVERSITAT JAUME I

Conclusions

# A synthetic tool for analysing adaptive workloads

**Authors**:
Iker Martín-Álvarez (martini@uji.es)
José I. Aliaga
María Isabel Castillo
Sergio Iserte
Rafael Mayo