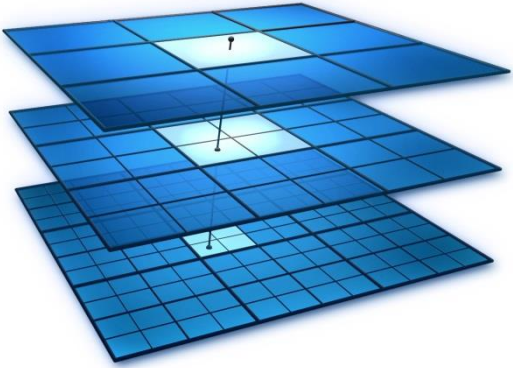
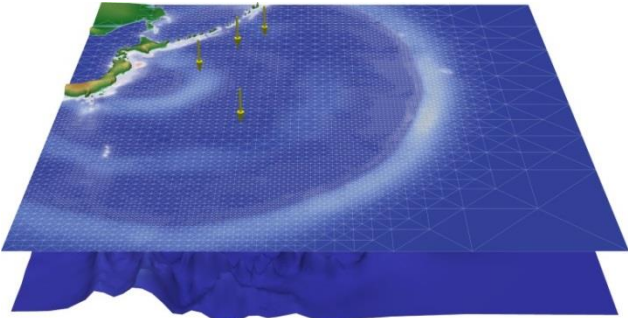


Infrastructure and API Extensions for Elastic Execution of MPI Applications

Isaías Comprés, Ao Mo-Hellenbrand,
Michael Gerndt, Hans-Joachim Bungartz

EuroMPI 2016, Edinburgh, Scotland

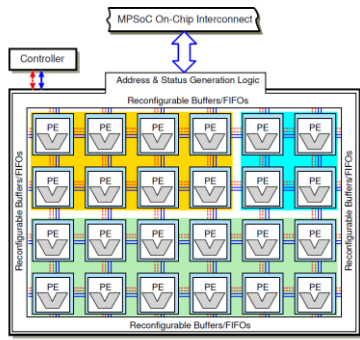
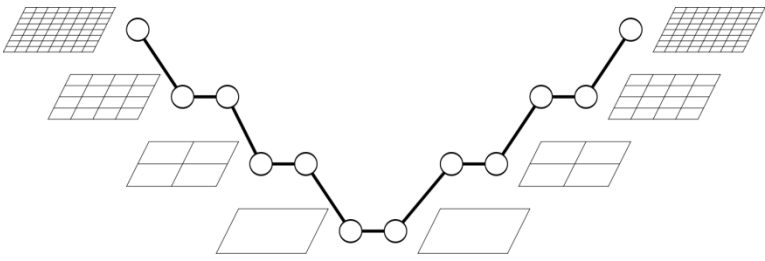
Invasive HPC applications



Invasive resource management for HPC



Invasive X10 applications on MPSoCs



Motivation

- State of the art scientific applications that utilize algorithms with evolving scalability properties
- The current assignment of fixed resources to these applications is suboptimal

Implementation

- **MPI extensions**
 - Extend the API with adaptive operations
- **MPI Library**
 - Based on MPICH 3.2
- **Resource Manager**
 - Based on SLURM 15.08

Proposed 4 new operations as an extension to the MPI standard:

`MPI_Init_adapt (...)`

- Initializes the library in adaptive mode

`MPI_Probe_adapt (...)`

- Probes the resource manager for adaptations

`MPI_Comm_adapt_begin (...)`

- Marks the beginning of an adaptation window
- Provides a set of helper communicators

`MPI_Comm_adapt_commit (...)`

- Marks the end of an adaptation window
- Sets adapted `MPI_COMM_WORLD`

Code Structure

```
MPI_Init_adapt (... , &status);
for (...) {
    MPI_Probe_adapt (&adapt, ...);
    if (adapt) {
        MPI_Comm_adapt_begin (...);
        // redistribution code
        MPI_Comm_adapt_commit (...);
    }
    // compute and MPI code
}
```

```
int MPI_Init_adapt (  
    int * argc ,  
    char *** argv ,  
    int * status  
);
```

status:

- New
- Joining

Code Structure

```
MPI_Init_adapt(..., &status);  
for (...) {  
    MPI_Probe_adapt(&adapt,...);  
    if(adapt) {  
        MPI_Comm_adapt_begin(...);  
        // redistribution code  
        MPI_Comm_adapt_commit(...);  
    }  
    // compute and MPI code  
}
```

```
int MPI_Probe_adapt (  
    int * operation,  
    int * status,  
    MPI_Info * info  
);
```

operation:

- Addition
- Reduction
- Combined
- Migration

status:

- Joining
- Staying
- Leaving

Code Structure

```
MPI_Init_adapt(..., &status);  
for (...) {  
    MPI_Probe_adapt(&adapt, ...);  
    if (adapt) {  
        MPI_Comm_adapt_begin(...);  
        // redistribution code  
        MPI_Comm_adapt_commit(...);  
    }  
    // compute and MPI code  
}
```

```
int MPI_Probe_adapt (
    int * operation,
    int * status,
    int * nfailed,
    int ** failed_ranks,
    MPI_Info * info
);
```

operation:

- Expansion
- Reduction
- Combined
- Migration
- **Fault**

status:

- New
- Joining
- Staying
- Leaving

Code Structure

```
MPI_Init_adapt(..., &status);
for (...) {
    MPI_Probe_adapt(&adapt, ...);
    if (adapt) {
        MPI_Comm_adapt_begin(...);
        // redistribution code
        MPI_Comm_adapt_commit(...);
    }
    // compute and MPI code
}
```

```
int MPI_Comm_adapt_begin (  
    int * intercomm,  
    int * future_comm_world  
);
```

intercomm:

- Equivalent to `MPI_Comm_spawn`
- Used to reach parents from the child group, and the children from the parent group

future_comm_world:

- Contains all **staying** parents plus children processes
- Parent processes that are **leaving** receive `MPI_COMM_NULL`

Code Structure

```
MPI_Init_adapt(..., &status);  
for (...) {  
    MPI_Probe_adapt(&adapt, ...);  
    if (adapt) {  
        MPI_Comm_adapt_begin(...);  
        // redistribution code  
        MPI_Comm_adapt_commit(...);  
    }  
    // compute and MPI code  
}
```



```
int MPI_Comm_adapt_commit ();
```

- `MPI_COMM_WORLD` is set to the `new_comm_world` communicator provided by the **`MPI_Comm_adapt_begin`** operation earlier.
- Leaving processes are required to terminate
 - In our current prototype the operation itself calls `exit()`
 - Our current applications clean up memory and file descriptors in the adaptation window, before commit

Code Structure

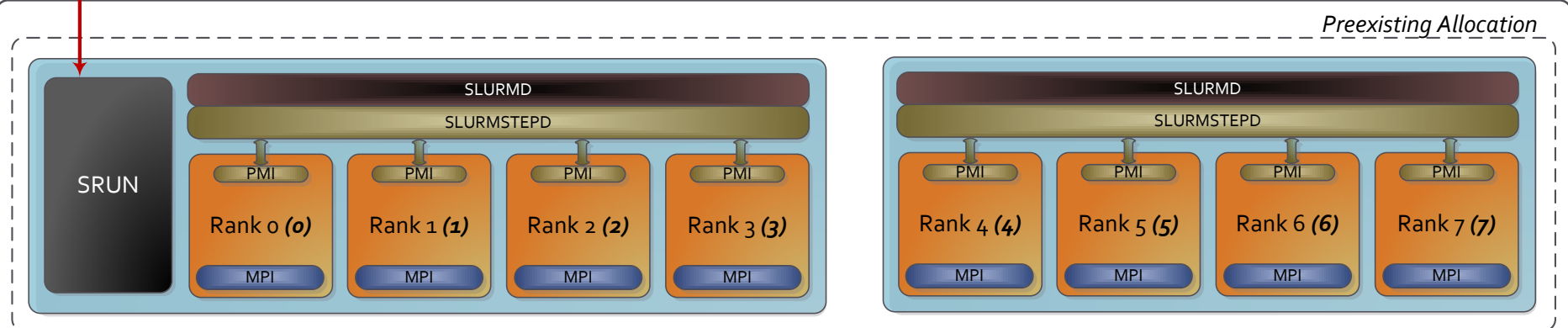
```
MPI_Init_adapt(..., &status);  
for (...) {  
    MPI_Probe_adapt(&adapt, ...);  
    if (adapt) {  
        MPI_Comm_adapt_begin(...);  
        // redistribution code  
        MPI_Comm_adapt_commit(...);  
    }  
    // compute and MPI code  
}
```

1: Reallocation Message

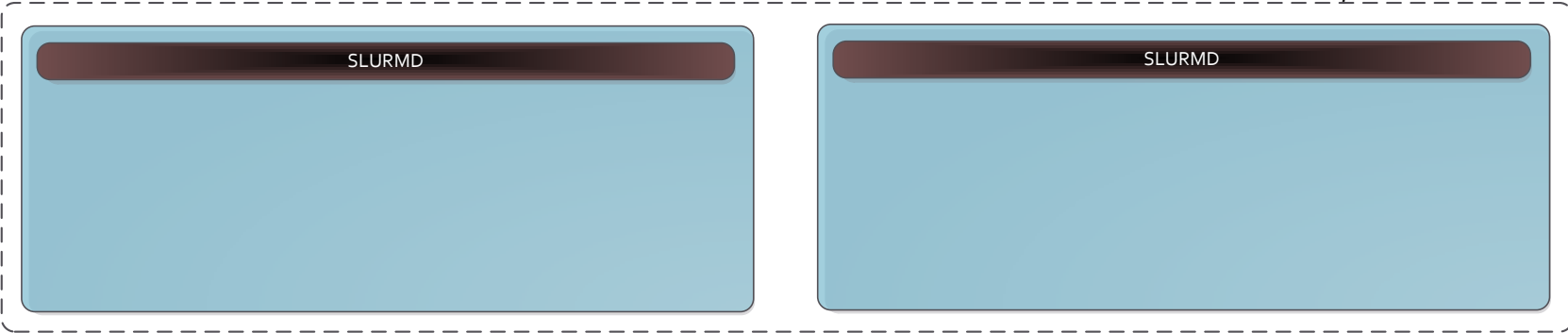


New Adapted Allocation

Preexisting Allocation



Expansion Allocation

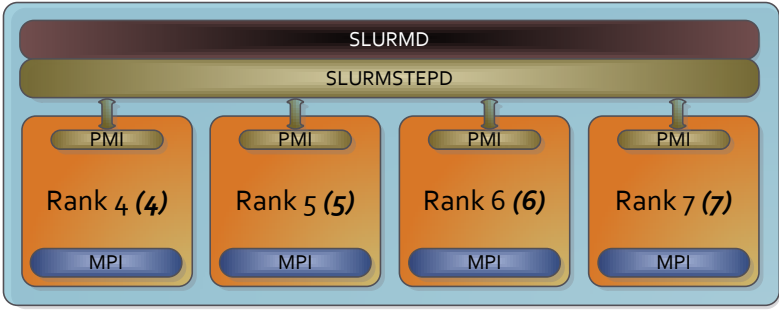
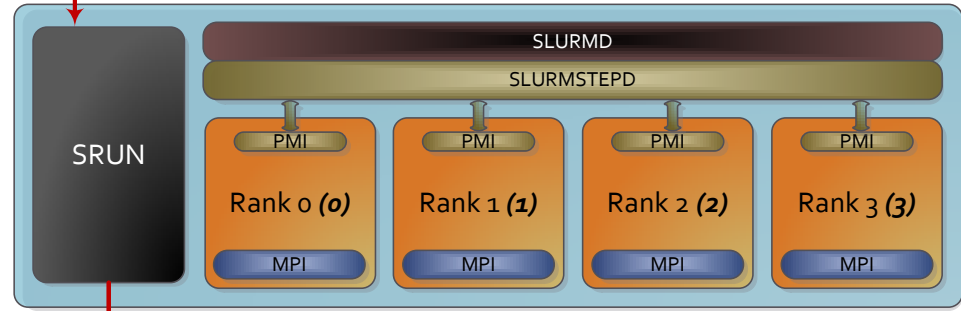


1: Reallocation Message



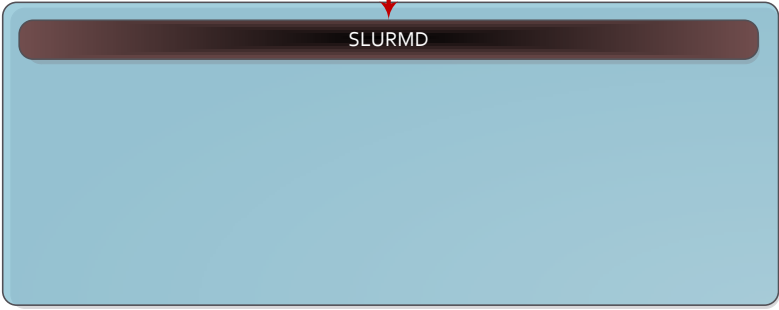
New Adapted Allocation

Preexisting Allocation



2: Create New Processes in Expansion Nodes

Expansion Allocation

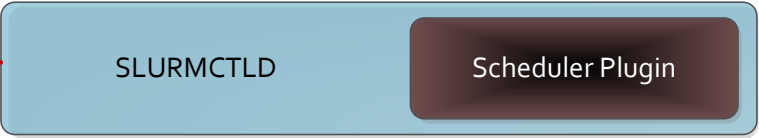


D3

Adaptation Step 3

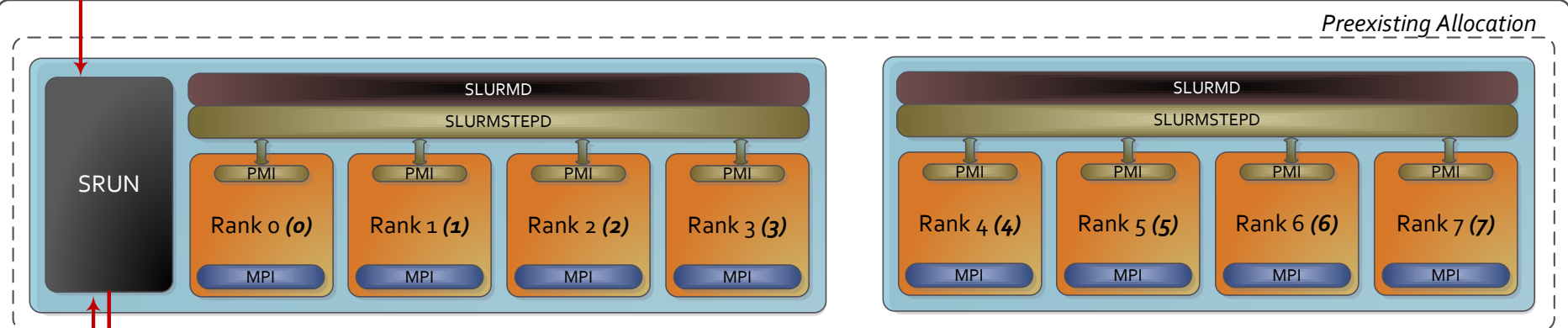


1: Reallocation Message



New Adapted Allocation

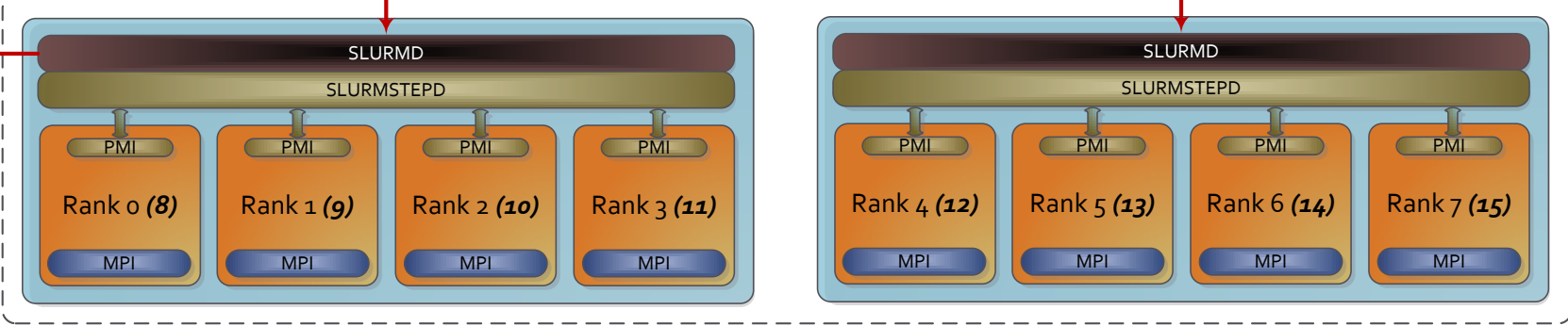
Preexisting Allocation



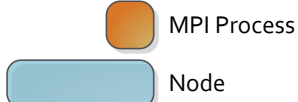
Expansion Allocation

3: New Processes Ready

2: Create New Processes in Expansion Nodes

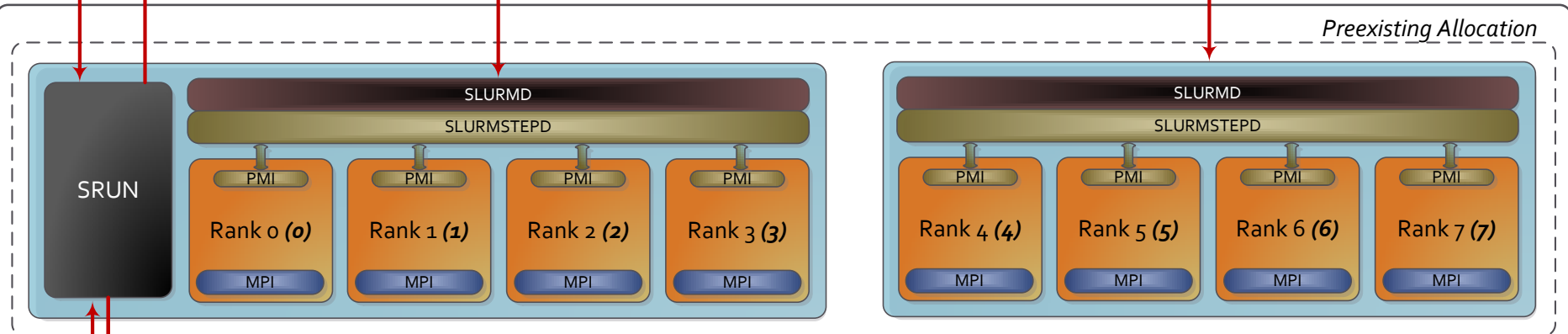


1: Reallocation Message



4: Notify Preexisting Processes

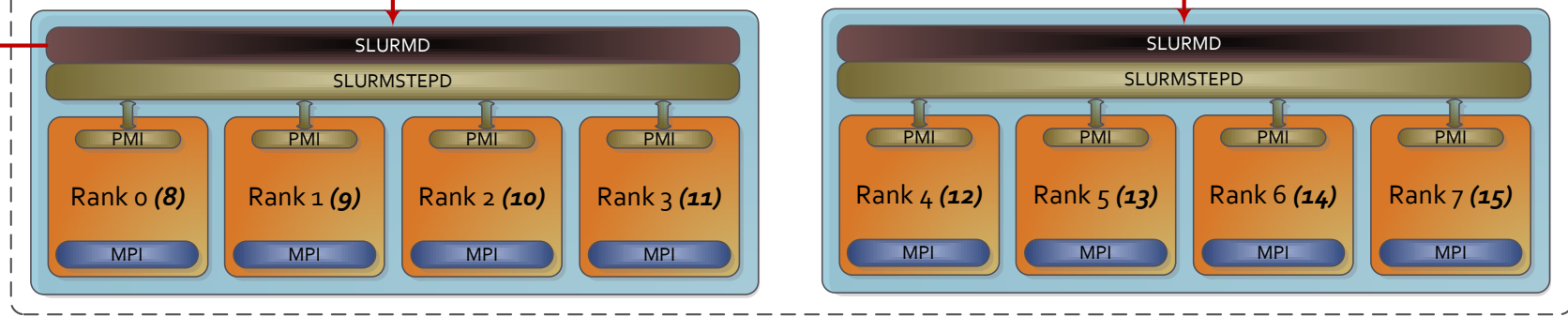
New Adapted Allocation
Preexisting Allocation



3: New Processes Ready

2: Create New Processes in Expansion Nodes

Expansion Allocation

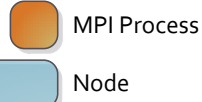


D3

Adaptation Step 5



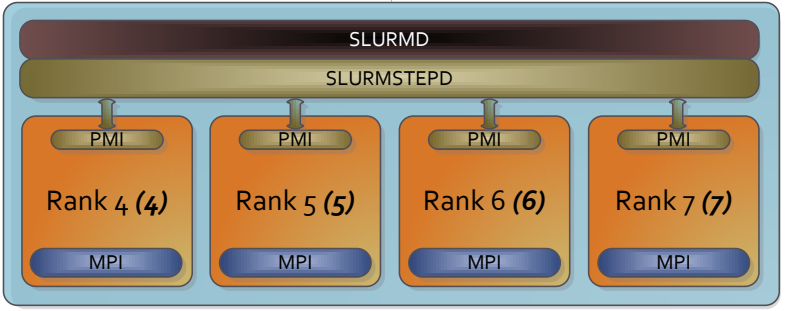
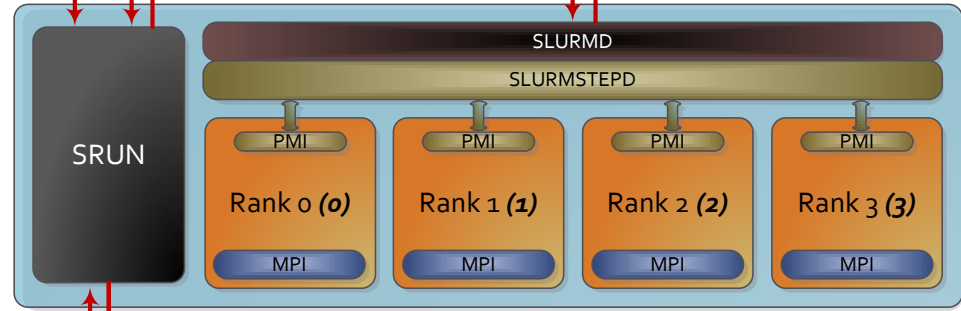
1: Reallocation Message



5: Adaptation Commit

4: Notify Preexisting Processes

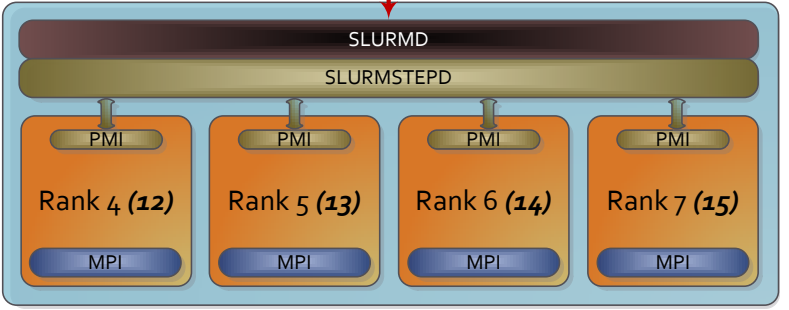
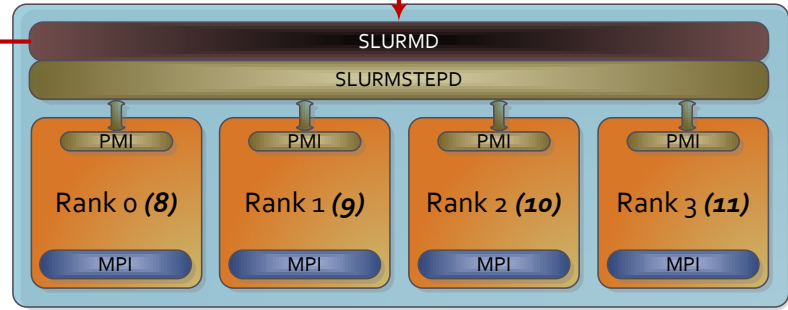
New Adapted Allocation
Preexisting Allocation

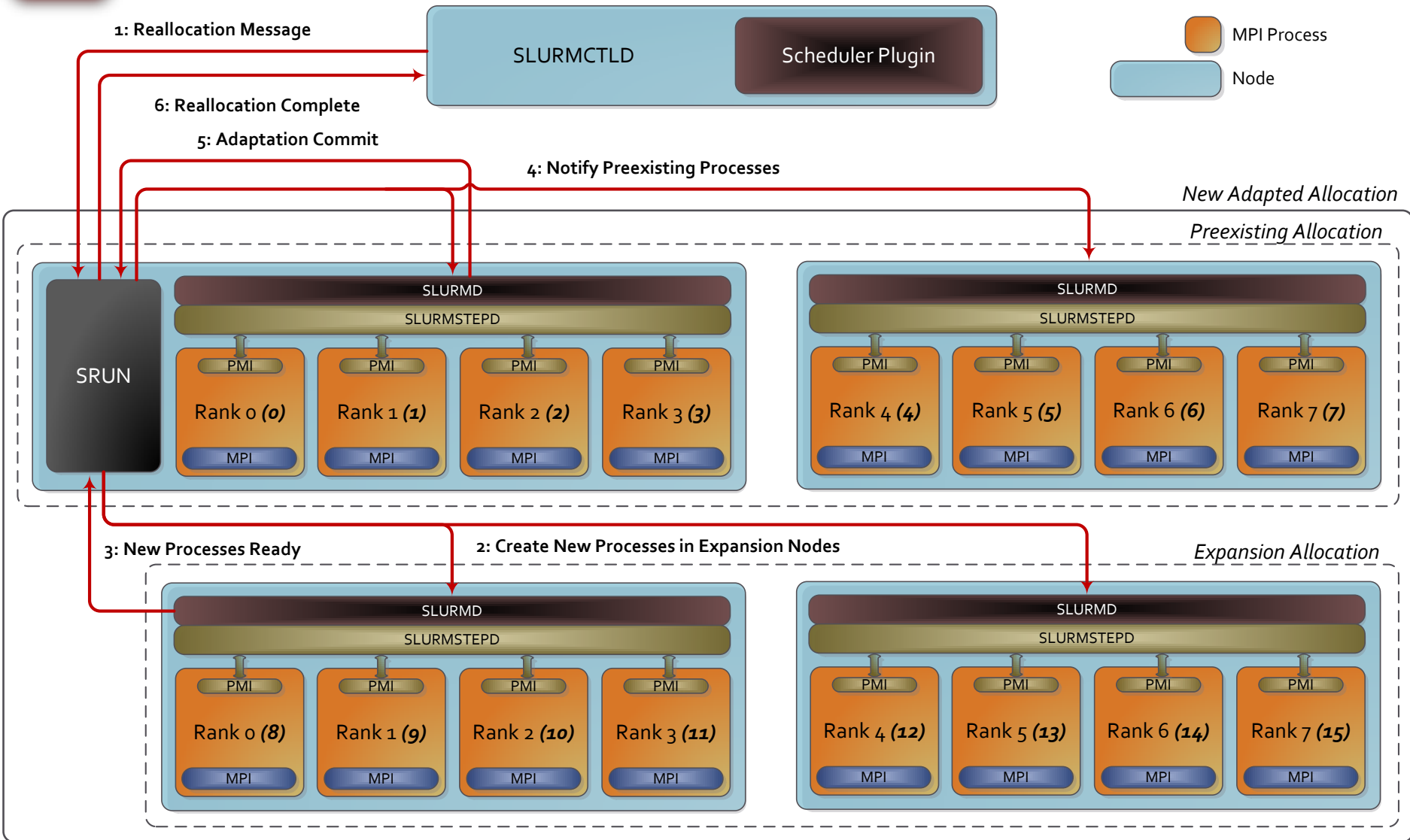


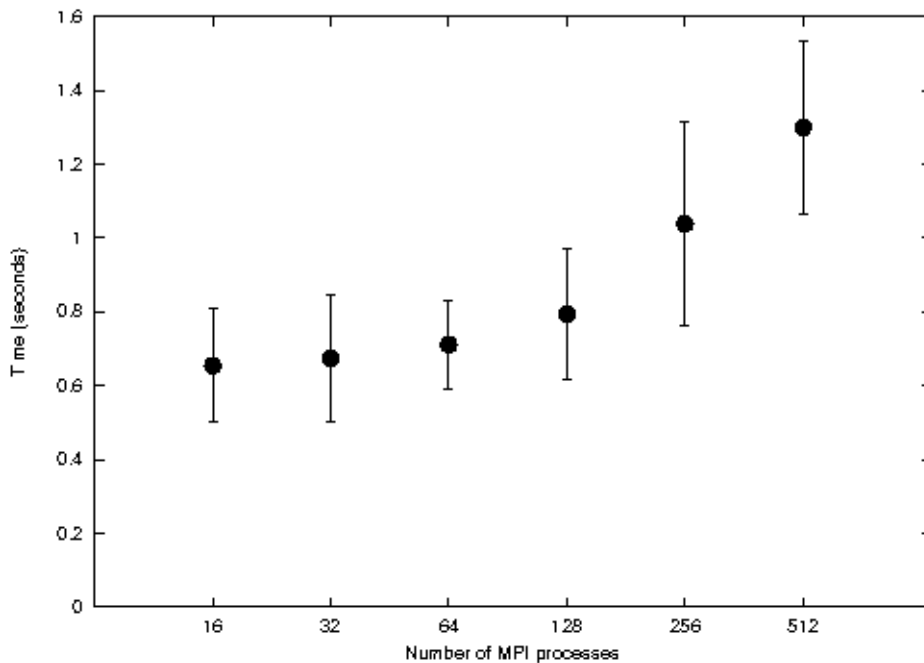
3: New Processes Ready

2: Create New Processes in Expansion Nodes

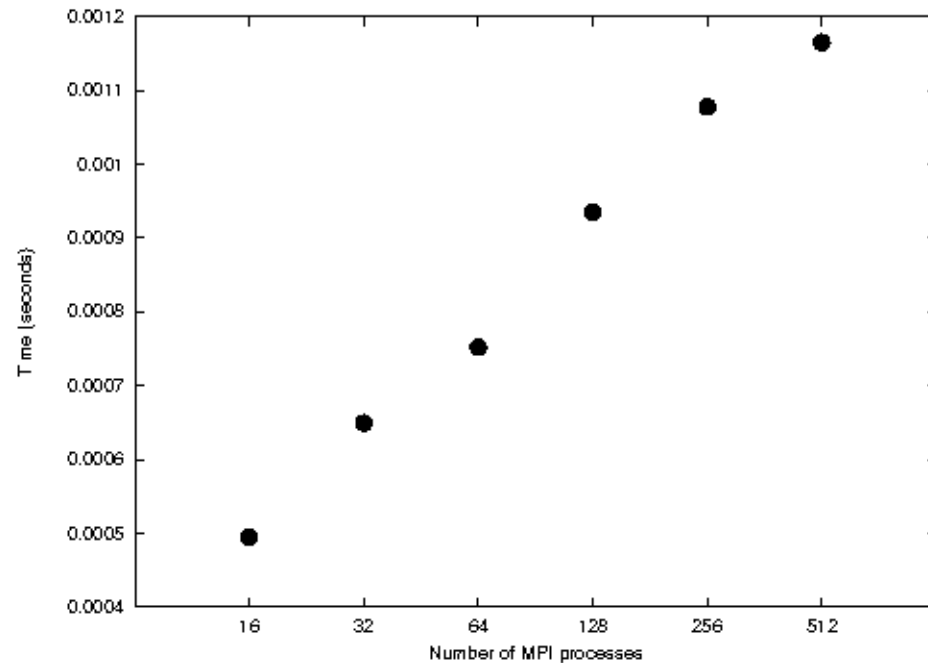
Expansion Allocation



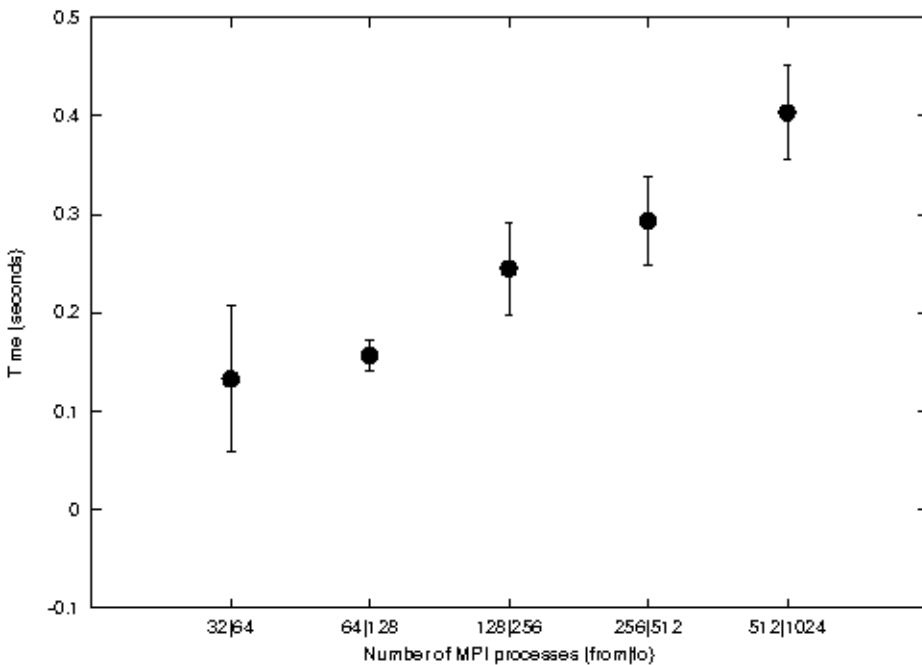




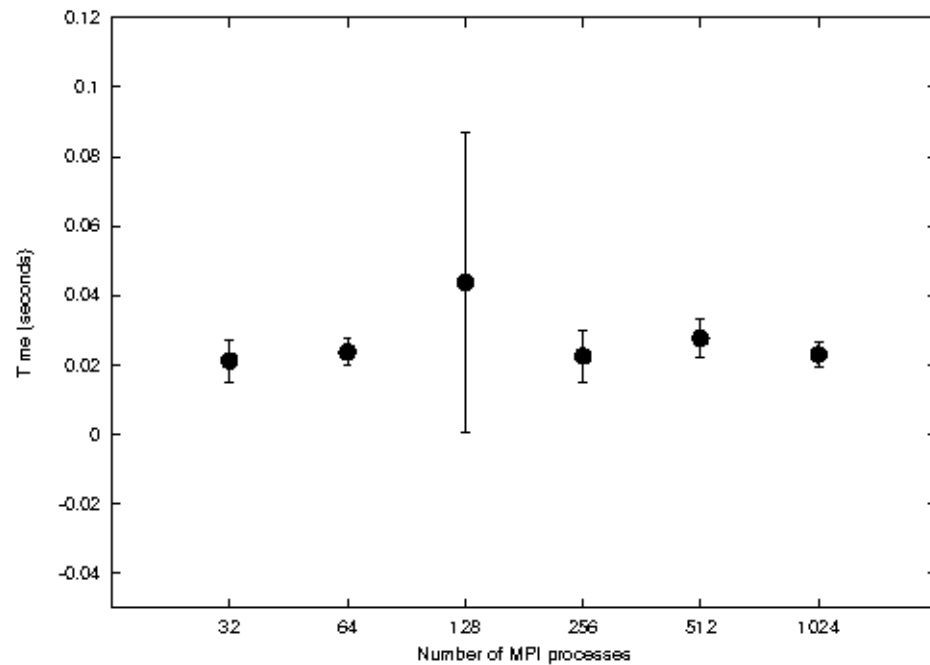
MPI_Init_adapt



MPI_Probe_adapt



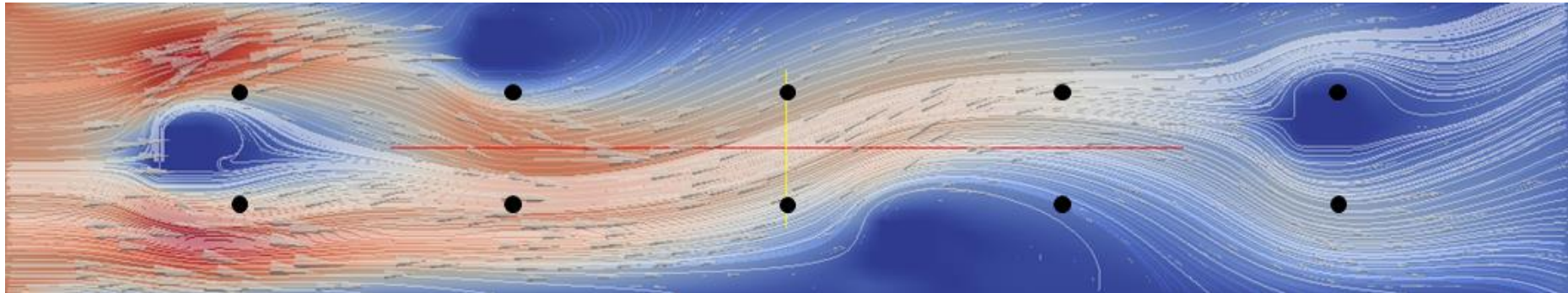
MPI_Comm_adapt_begin

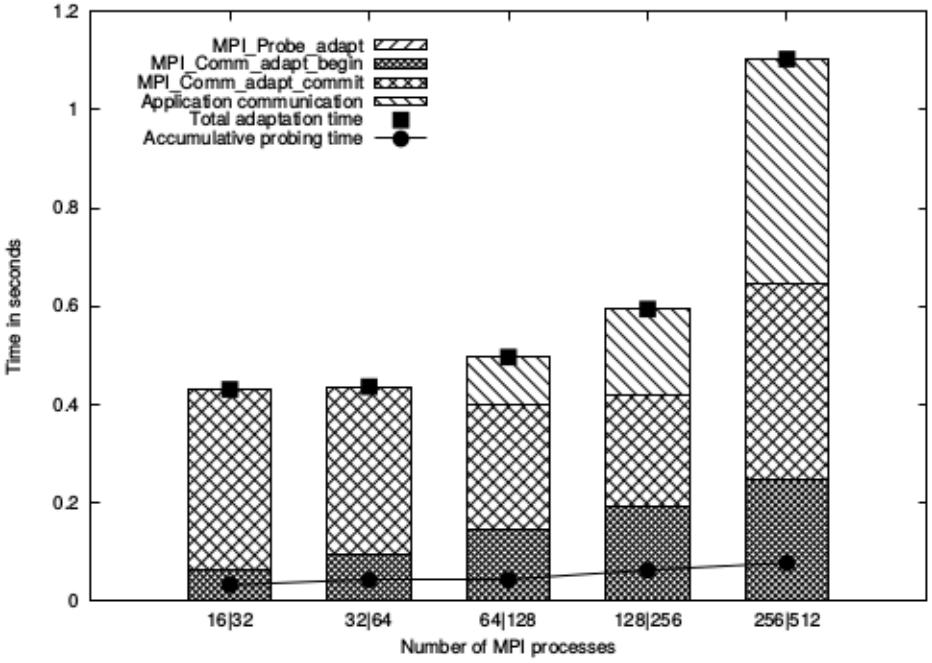


MPI_Comm_adapt_commit

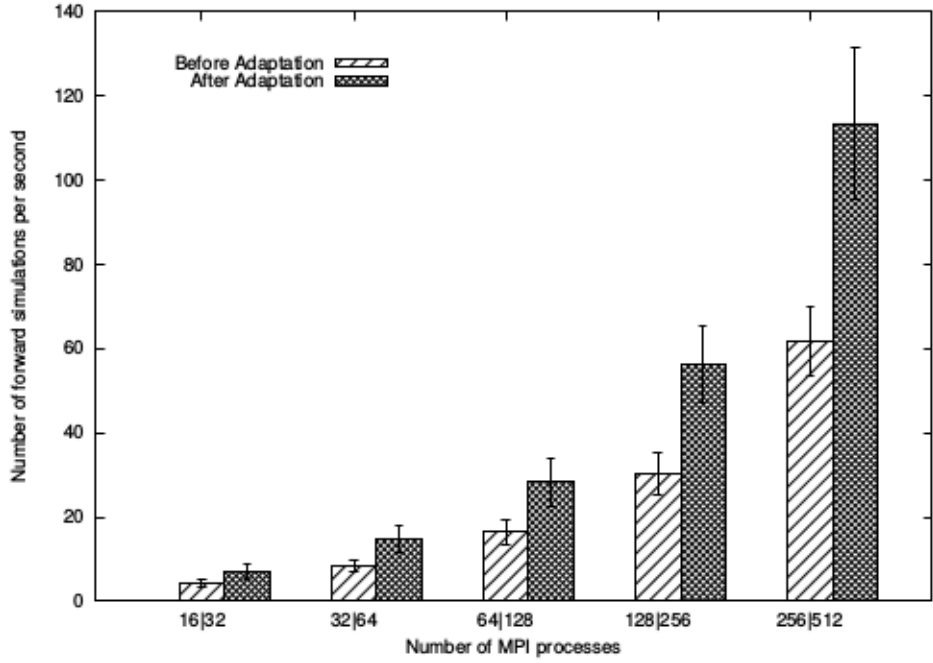
Elastic Surrogate Model Construction for a Statistical Inverse Problem

- Locates the number of obstacles in a fluid channel
- 2D version as first elastic conversion
- Fluid simulation as input, instead of real physical flow
 - Quick setup of various experimental scenarios
 - Easy verification of the method's success
- Outputs the predicted obstacles location

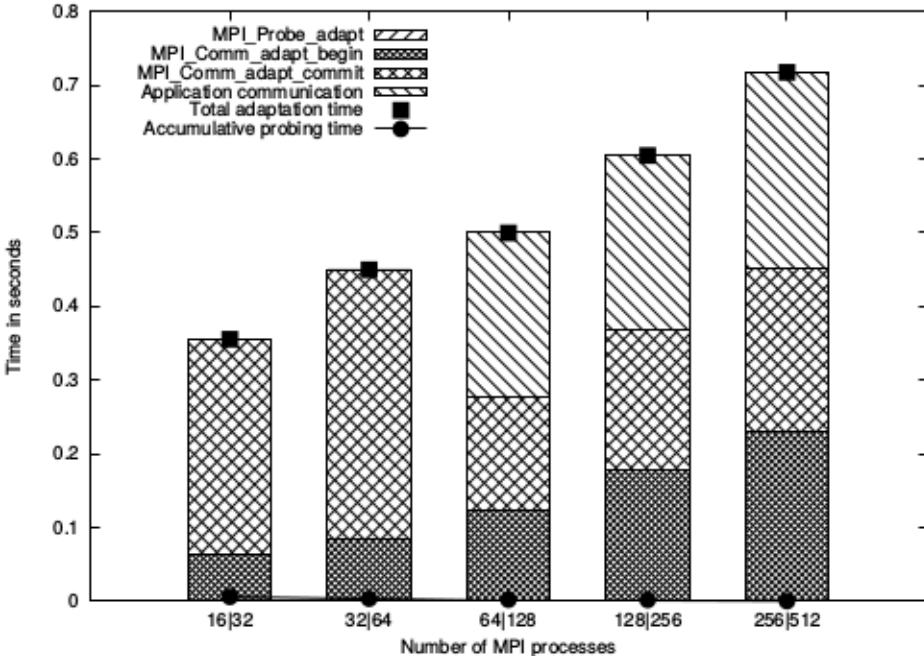




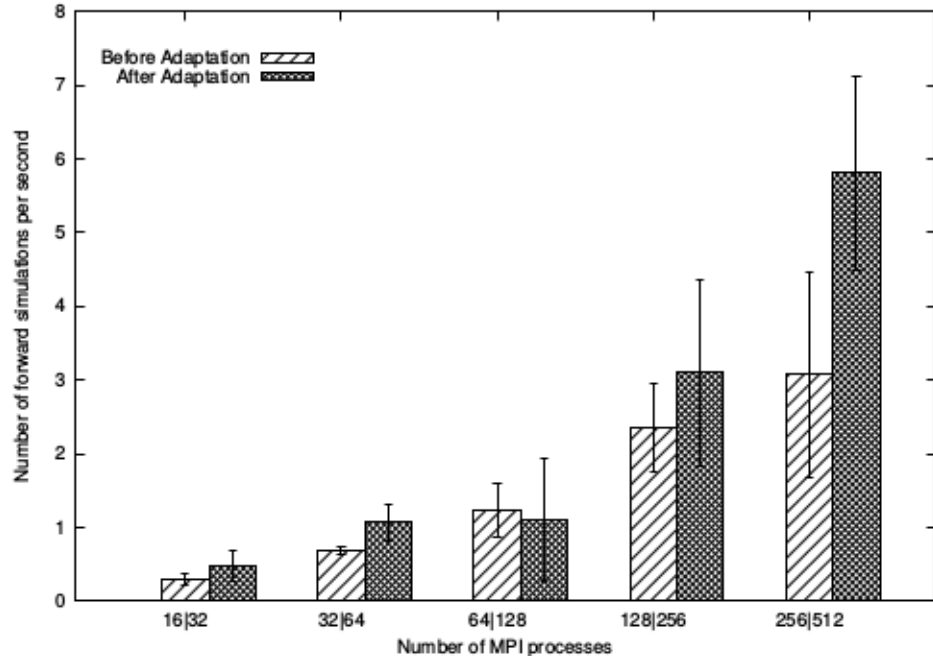
Adaptation overhead (small)



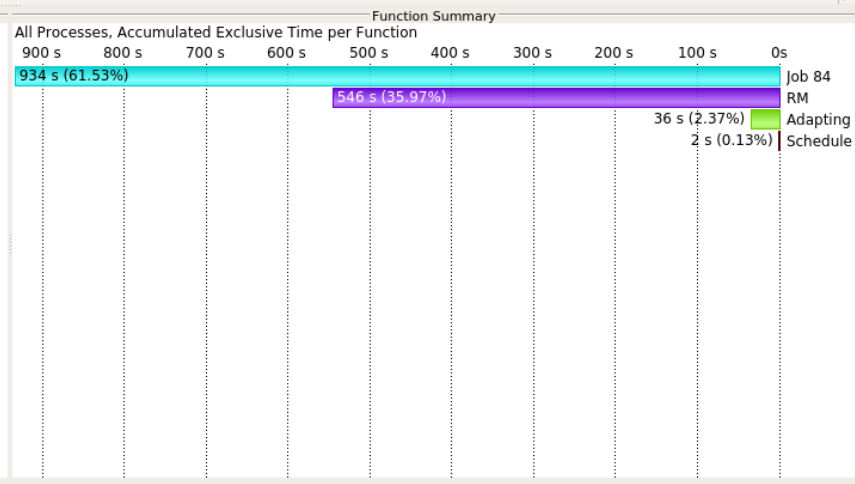
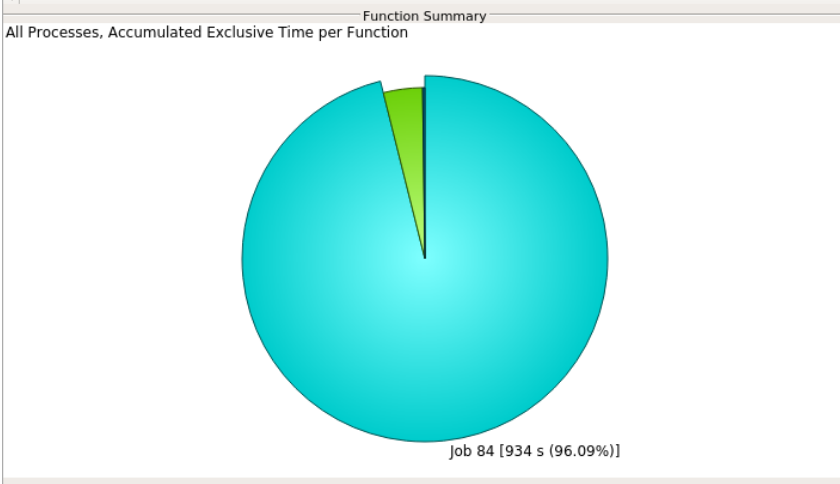
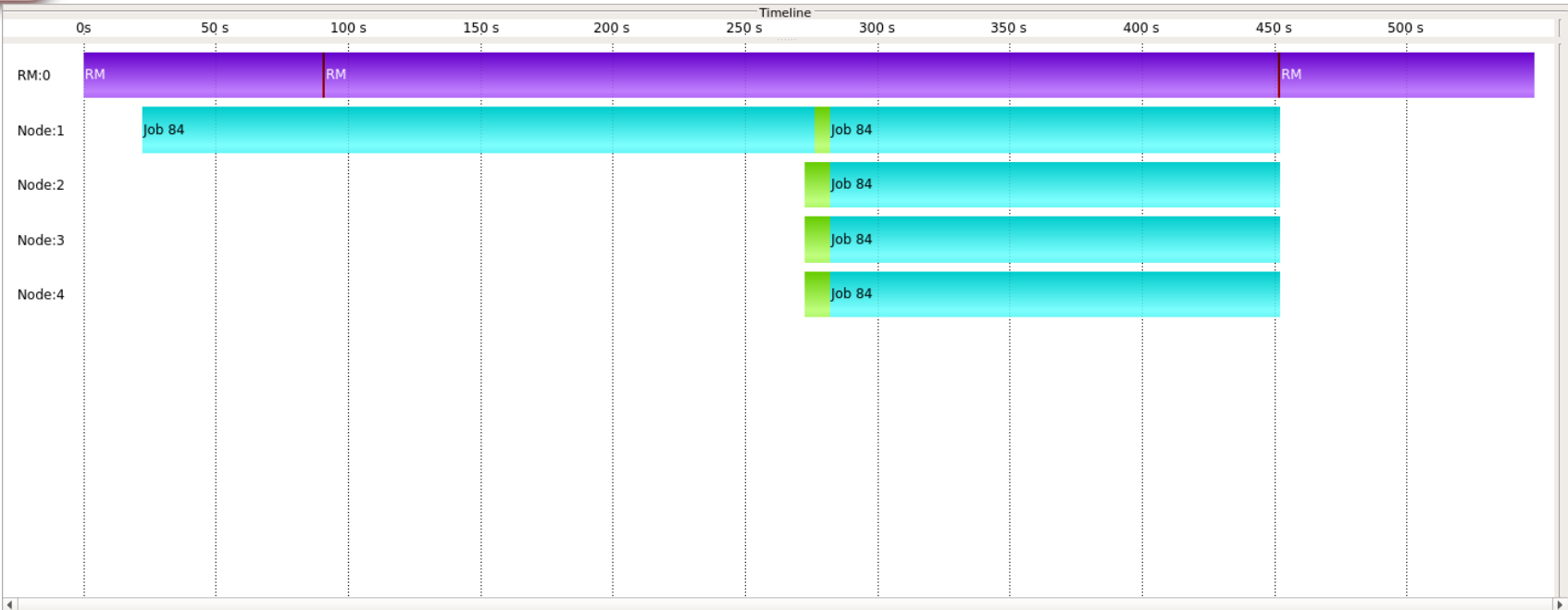
Performance change (small)

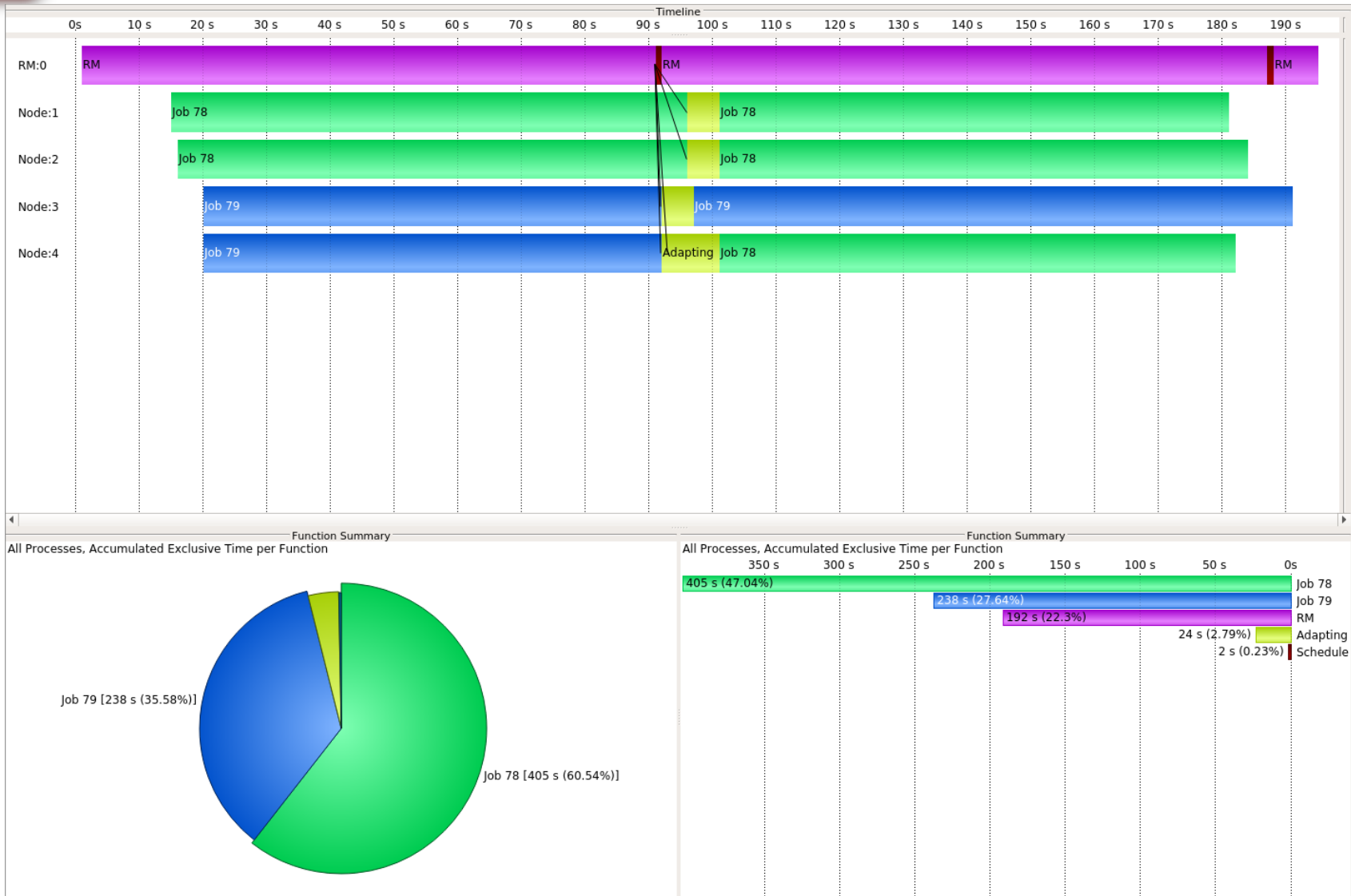


Adaptation overhead (large)



Performance change (large)





Proposed an extension to MPI that consists of 4 operations for:

- Initialization
- Probing for adaptations, and
- The creation of adaptation windows

Supporting these operations required changes to the resource manager and the PMI, in addition to the MPI library.

Our prototype implementation based on MPICH and SLURM provides:

- Satisfactory performance for our current use cases
- Latency hiding properties that minimize waiting times in preexisting application processes

Future work:

- Scheduling and tuning based on performance modeling
- Potential move to new resource manager (e.g. Flux) and PMI implementations (e.g. PMIx)