

Invading instead of Wasting HPC resources

Michael Gerndt
Technische Universität München

10th VI-HPS Anniversary, Frankfurt



Invasive Resource Management

Dynamic

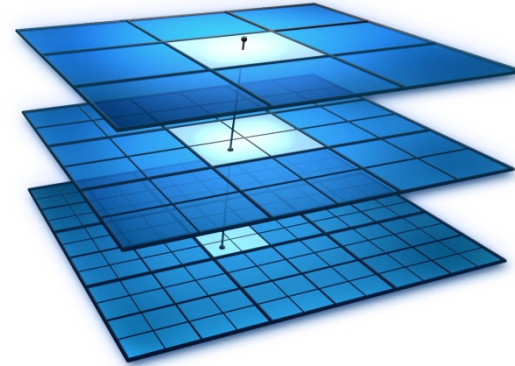
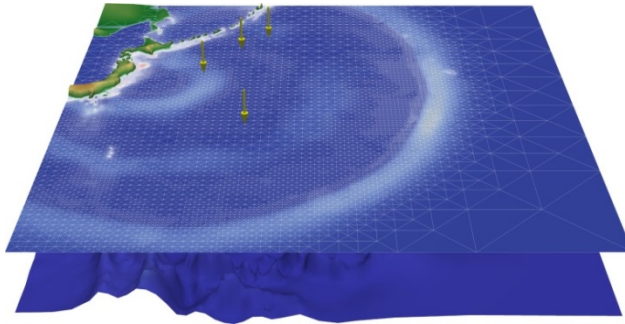
- Applications
- Capabilities
- Resources

Static Resource Management

Scenarios:

- Phase-based applications
- Adaptive applications
- Jobs with different input data set
- Analytics applications with data bursts
- Coupled applications with dynamically varying resource requirements
- IO intensive phases require network and IO bandwidth
- Node failures
- Urgent computing
- Power stability
- Increased scheduling opportunities
- ...

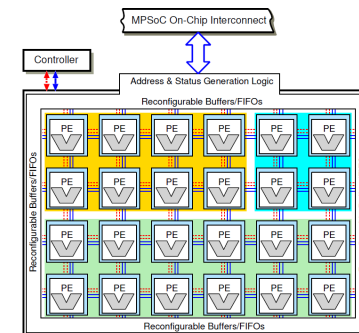
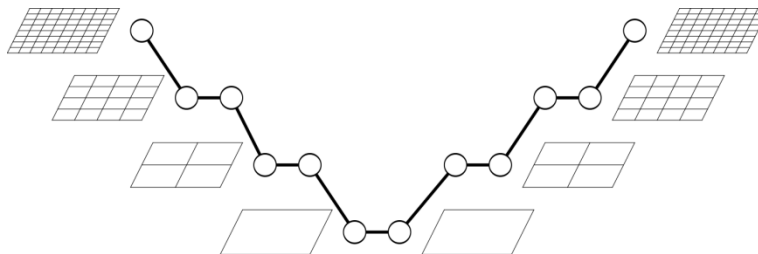
Invasive HPC applications



Invasive resource management for HPC



Invasive X10 applications on MPSoCs



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 inter communicator and new communicator

`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  
}
```

1: Reallocation Message

SLURMCTLD

Scheduler Plugin

 MPI Process Node

New Adapted Allocation

Preexisting Allocation

SRUN

SLURMD

SLURMSTEPD

PMI

Rank 0 (0)

MPI

PMI

Rank 1 (1)

MPI

PMI

Rank 2 (2)

MPI

PMI

Rank 3 (3)

MPI

SLURMD

SLURMSTEPD

PMI

Rank 4 (4)

MPI

PMI

Rank 5 (5)

MPI

PMI

Rank 6 (6)

MPI

PMI

Rank 7 (7)

MPI

SLURMD

SLURMD

Expansion Allocation

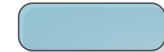
1: Reallocation Message

SLURMCTLD

Scheduler Plugin



MPI Process



Node

New Adapted Allocation

Preexisting Allocation

SRUN

SLURMD

SLURMSTEPD

PMI

Rank 0 (0)

MPI

PMI

Rank 1 (1)

MPI

PMI

Rank 2 (2)

MPI

PMI

Rank 3 (3)

MPI

SLURMD

SLURMSTEPD

PMI

Rank 4 (4)

MPI

PMI

Rank 5 (5)

MPI

PMI

Rank 6 (6)

MPI

PMI

Rank 7 (7)

MPI

2: Create New Processes in Expansion Nodes

SLURMD

SLURMD

Expansion Allocation

1: Reallocation Message

SLURMCTLD

Scheduler Plugin

 MPI Process Node

New Adapted Allocation

Preexisting Allocation

SRUN

SLURMD

SLURMSTEPD

PMI

Rank 0 (0)

PMI

Rank 1 (1)

PMI

Rank 2 (2)

PMI

Rank 3 (3)

MPI

MPI

MPI

MPI

SLURMD

SLURMSTEPD

PMI

Rank 4 (4)

PMI

Rank 5 (5)

PMI

Rank 6 (6)

PMI

Rank 7 (7)

MPI

MPI

MPI

MPI

3: New Processes Ready

2: Create New Processes in Expansion Nodes

Expansion Allocation

SLURMD

SLURMSTEPD

PMI

Rank 0 (8)

PMI

Rank 1 (9)

PMI

Rank 2 (10)

PMI

Rank 3 (11)

MPI

MPI

MPI

MPI

SLURMD

SLURMSTEPD

PMI

Rank 4 (12)

PMI

Rank 5 (13)

PMI

Rank 6 (14)

PMI

Rank 7 (15)

MPI

MPI

MPI

MPI

1: Reallocation Message

SLURMCTLD

Scheduler Plugin

 MPI Process Node

4: Notify Preexisting Processes

New Adapted Allocation

Preexisting Allocation

SRUN

SLURMD

SLURMSTEPD

PMI

Rank 0 (0)

MPI

PMI

Rank 1 (1)

MPI

PMI

Rank 2 (2)

MPI

PMI

Rank 3 (3)

MPI

SLURMD

SLURMSTEPD

PMI

Rank 4 (4)

MPI

PMI

Rank 5 (5)

MPI

PMI

Rank 6 (6)

MPI

PMI

Rank 7 (7)

MPI

3: New Processes Ready

2: Create New Processes in Expansion Nodes

Expansion Allocation

SLURMD

SLURMSTEPD

PMI

Rank 0 (8)

MPI

PMI

Rank 1 (9)

MPI

PMI

Rank 2 (10)

MPI

PMI

Rank 3 (11)

MPI

SLURMD

SLURMSTEPD

PMI

Rank 4 (12)

MPI

PMI

Rank 5 (13)

MPI

PMI

Rank 6 (14)

MPI

PMI

Rank 7 (15)

MPI

1: Reallocation Message

SLURMCTLD

Scheduler Plugin

 MPI Process Node

5: Adaptation Commit

4: Notify Preexisting Processes

New Adapted Allocation

Preexisting Allocation

SRUN

SLURMD
SLURMSTEPD

PMI

Rank 0 (0)

MPI

PMI

Rank 1 (1)

MPI

PMI

Rank 2 (2)

MPI

PMI

Rank 3 (3)

MPI

SLURMD
SLURMSTEPD

PMI

Rank 4 (4)

MPI

PMI

Rank 5 (5)

MPI

PMI

Rank 6 (6)

MPI

PMI

Rank 7 (7)

MPI

3: New Processes Ready

2: Create New Processes in Expansion Nodes

Expansion Allocation

SLURMD
SLURMSTEPD

PMI

Rank 0 (8)

MPI

PMI

Rank 1 (9)

MPI

PMI

Rank 2 (10)

MPI

PMI

Rank 3 (11)

MPI

SLURMD
SLURMSTEPD

PMI

Rank 4 (12)

MPI

PMI

Rank 5 (13)

MPI

PMI

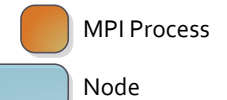
Rank 6 (14)

MPI

PMI

Rank 7 (15)

MPI



1: Reallocation Message

SLURMCTLD

Scheduler Plugin

6: Reallocation Complete

5: Adaptation Commit

4: Notify Preexisting Processes

New Adapted Allocation

Preexisting Allocation

SRUN

SLURMD
SLURMSTEPD

PMI

Rank 0 (0)

MPI

PMI

Rank 1 (1)

MPI

PMI

Rank 2 (2)

MPI

PMI

Rank 3 (3)

MPI

SLURMD
SLURMSTEPD

PMI

Rank 4 (4)

MPI

PMI

Rank 5 (5)

MPI

PMI

Rank 6 (6)

MPI

PMI

Rank 7 (7)

MPI

3: New Processes Ready

2: Create New Processes in Expansion Nodes

Expansion Allocation

SLURMD
SLURMSTEPD

PMI

Rank 0 (8)

MPI

PMI

Rank 1 (9)

MPI

PMI

Rank 2 (10)

MPI

PMI

Rank 3 (11)

MPI

SLURMD
SLURMSTEPD

PMI

Rank 4 (12)

MPI

PMI

Rank 5 (13)

MPI

PMI

Rank 6 (14)

MPI

PMI

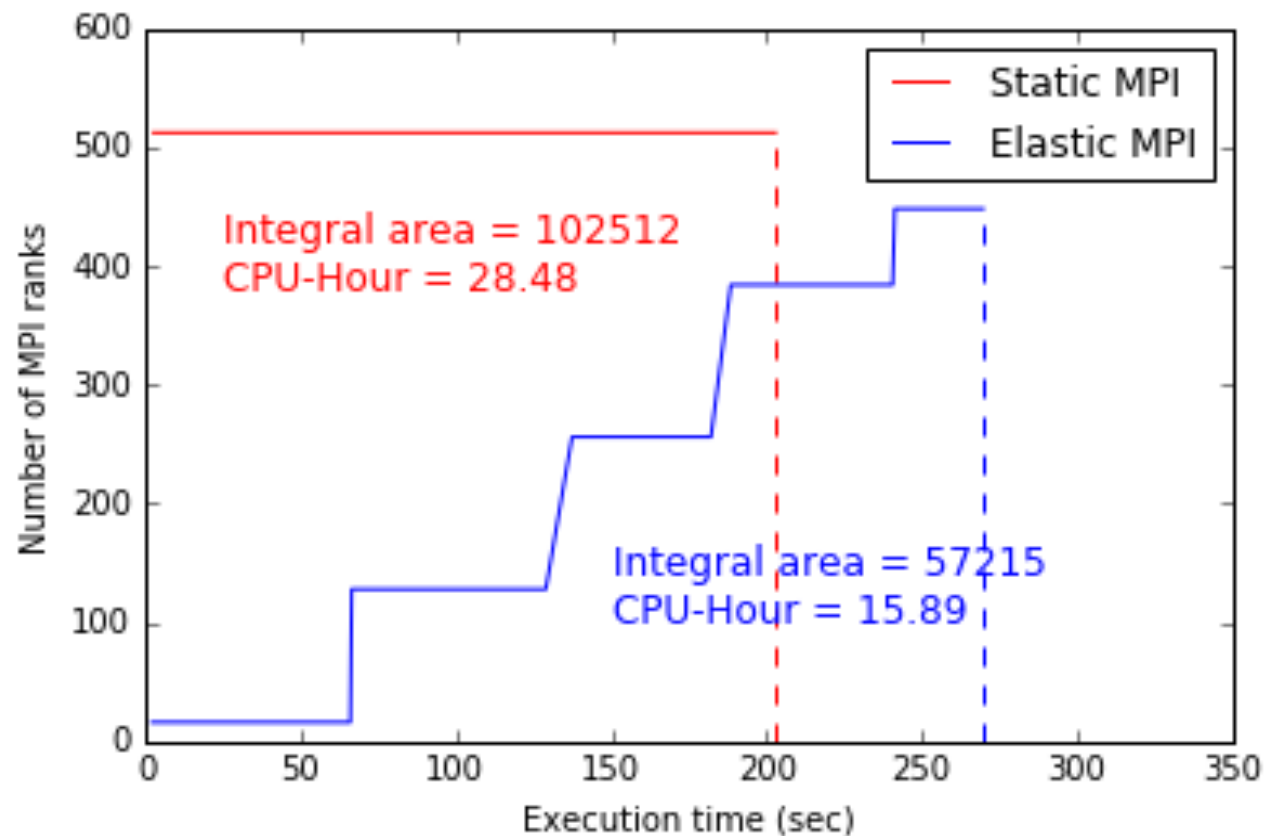
Rank 7 (15)

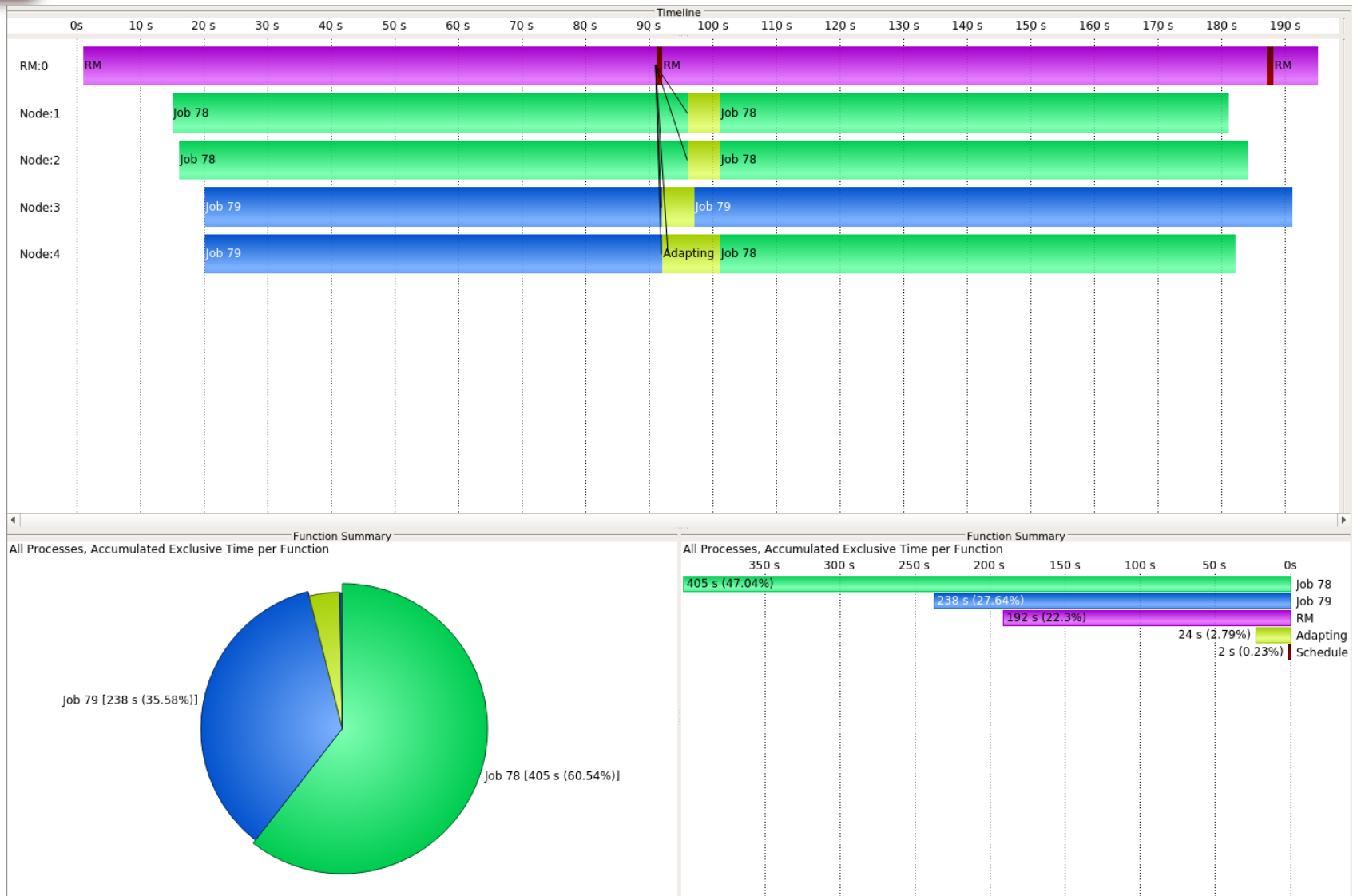
MPI

SuperMUC @ LRZ

- Allocation of set of SuperMUC nodes via batch job
- Management of the nodes via separate SLURM instance
- Distribution of resource management into
 - SLURM Scheduler
 - Selection and scheduling of invasive jobs
 - Based on resource offer
 - Invasic Scheduler
 - Invasive resource management
- Submission of new invasive jobs through sbatch command

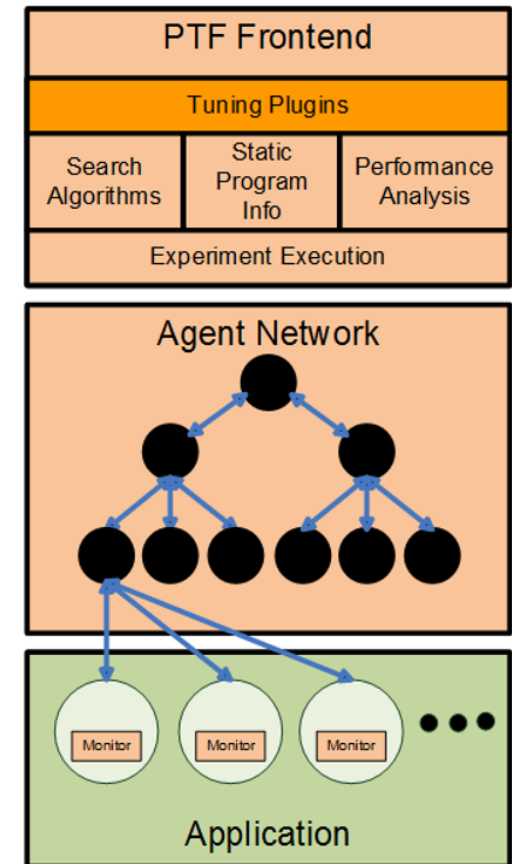
Tsunami simulations (Mo-Hellenbrand, Bungartz)





VI-HPS, AutoTune, READEX

- Periscope Frontend
 - Controls the analysis and tuning process
 - Performs a sequence of experiments
 - While the application is executing
 - Based on application phases
 - Automatically starting/restarting the application if required
- Agent Network for scalability
 - Leave agents responsible for a subset of the MPI processes
 - Intermediate agents aggregate performance properties
- Online Access to the monitoring system
 - Configuration of measurements and tuning actions
 - Retrieval of performance data



Thank You