

Performance Models for Dynamic Tuning of Parallel Applications on Computational Grids

Genaro Costa, Josep Jorba, Anna Morajko,
Tomas Margalef and Emilio Luque



Cluster 2008



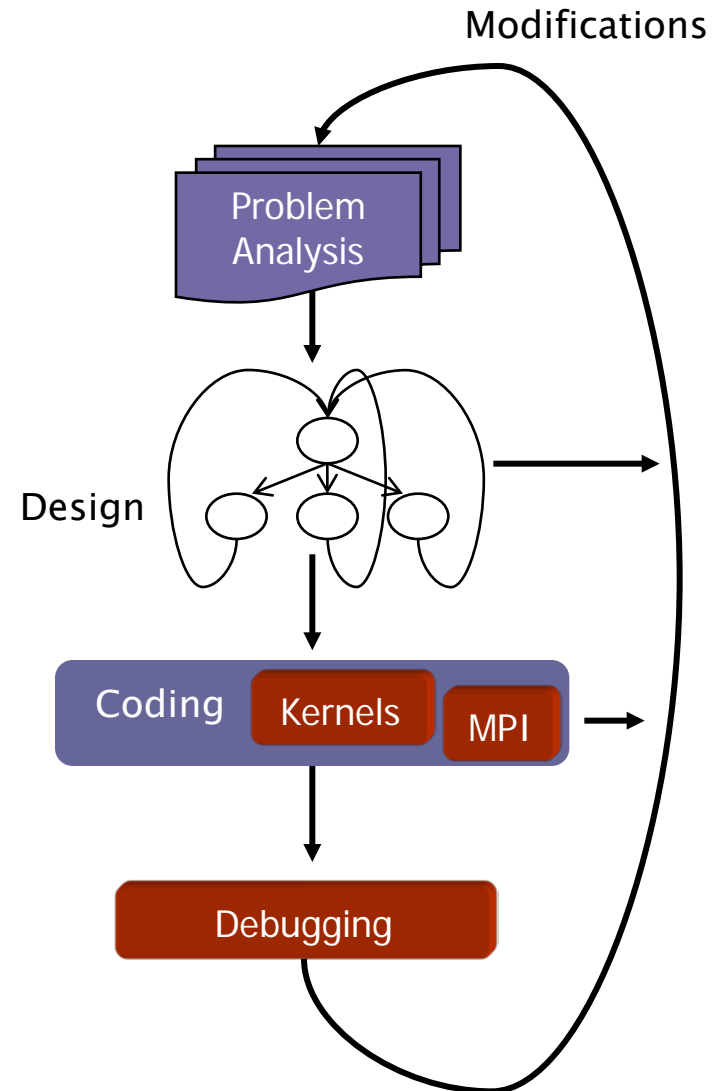
Content

1. Motivation
2. Performance Model for Dynamic Tuning
3. Tuning Heuristics
4. Experimentation Results
5. Conclusions

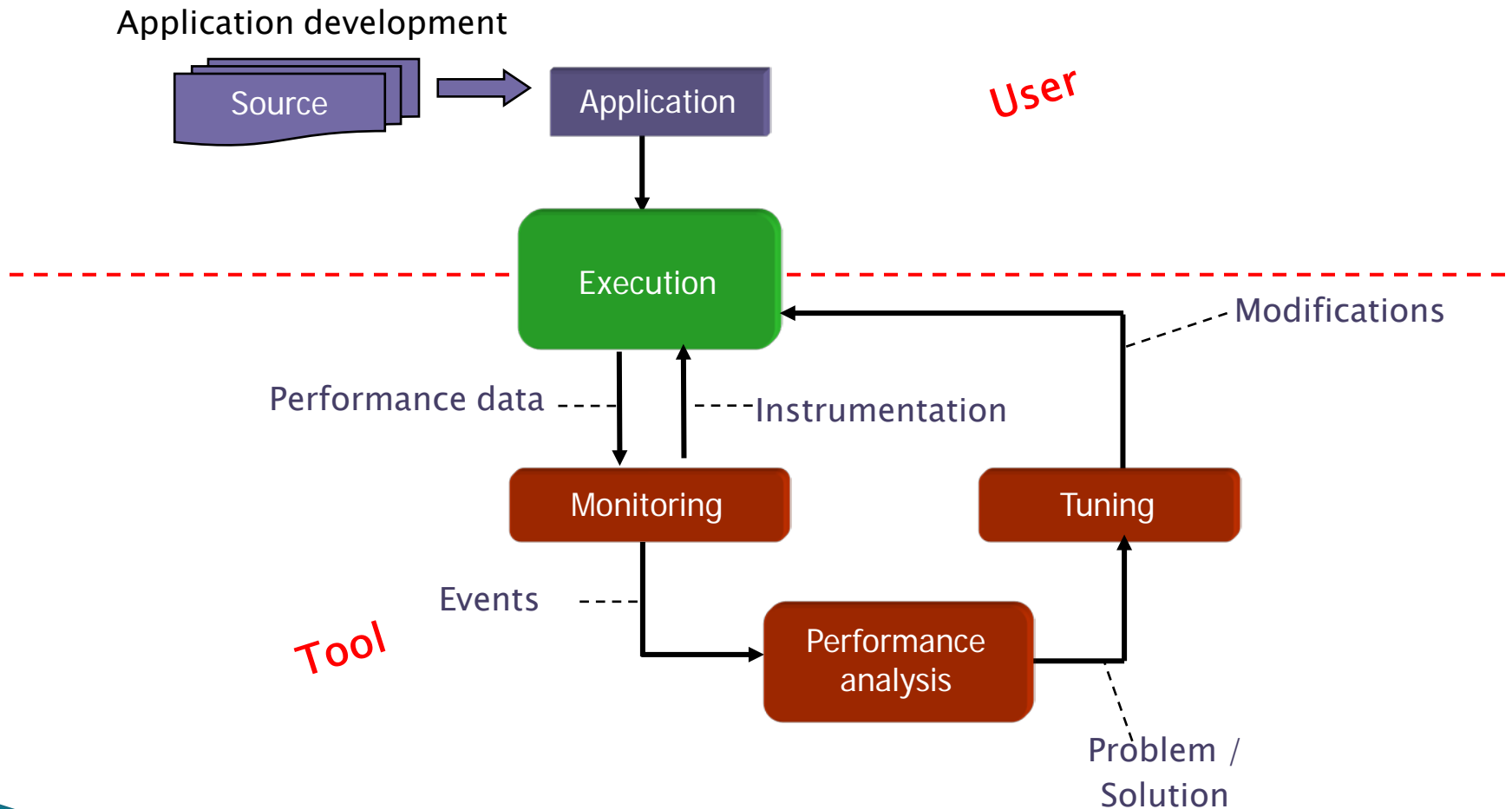


Motivation

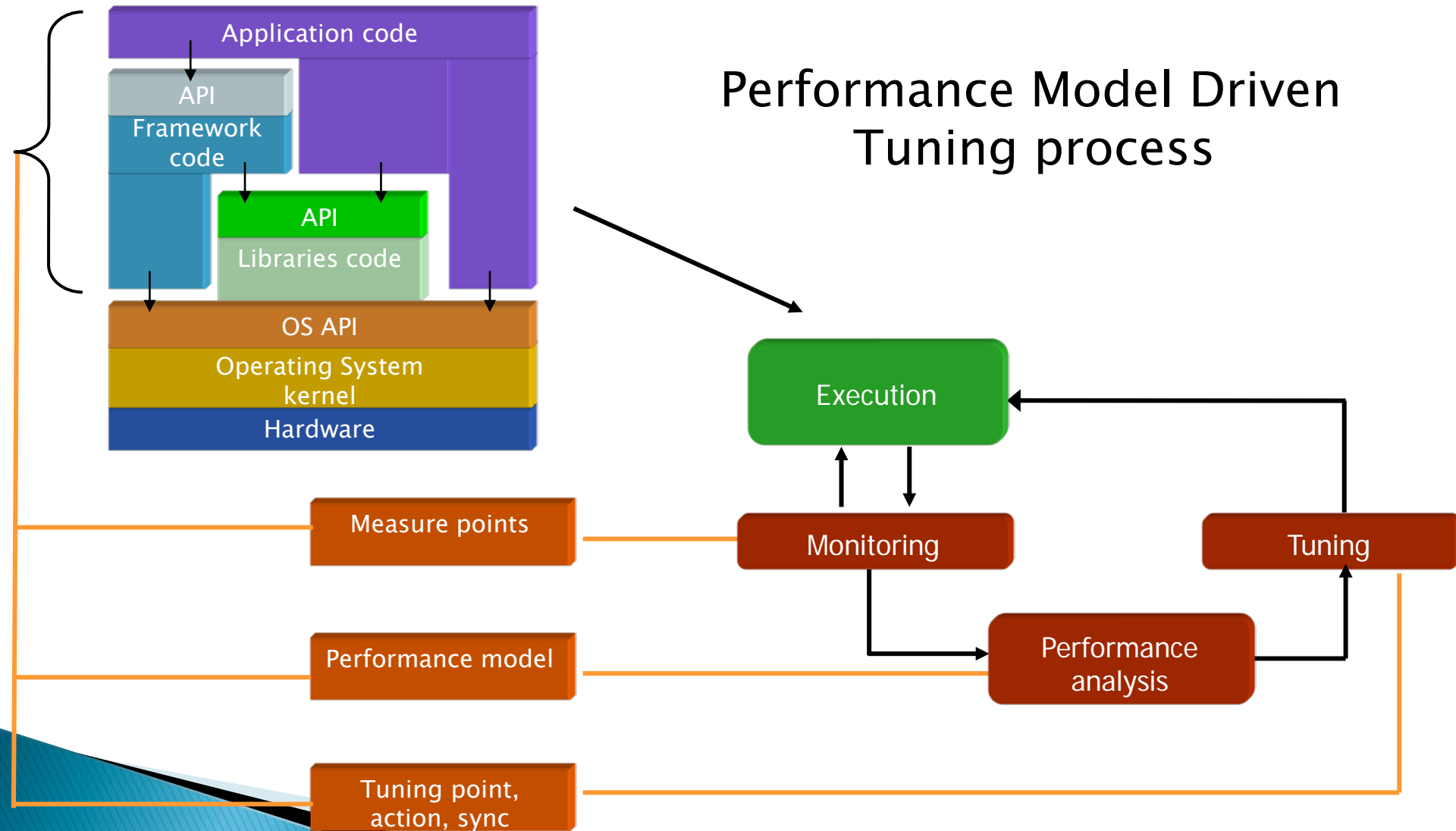
- ▶ Performance is a main issue in parallel application development.
- ▶ Runtime variation of system characteristics harms application performance.
- ▶ Computational Grids are heterogeneous and dynamic by nature.
- ▶ Dynamic tuning may help application increasing its execution performance indexes.



Automatic Performance Tuning



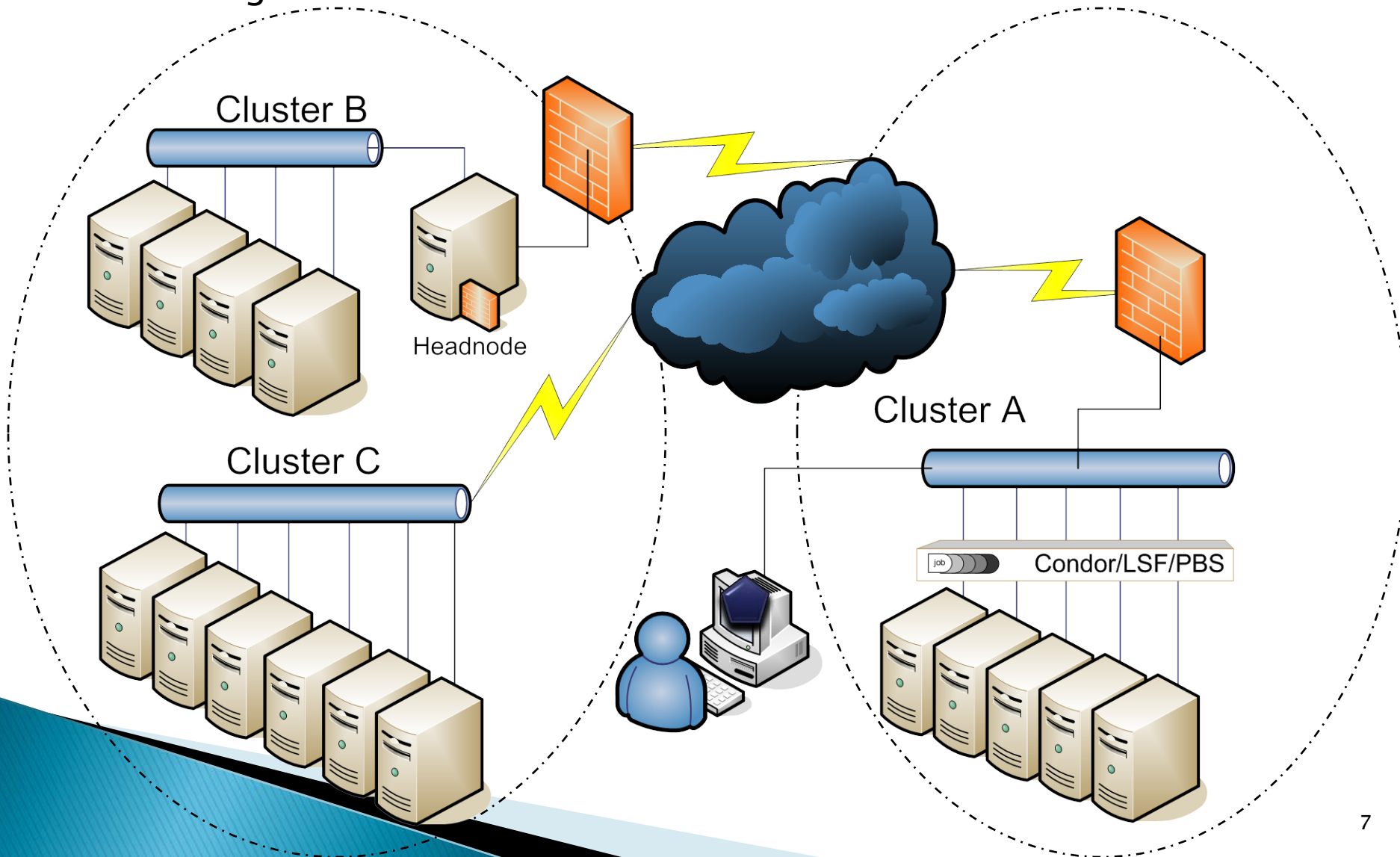
Dynamic Performance Tuning



Computational Grids

Organization A

Organization B



Content

1. Motivation

2. Performance Model for Dynamic Tuning

3. Tuning Heuristics

4. Experimentation Results

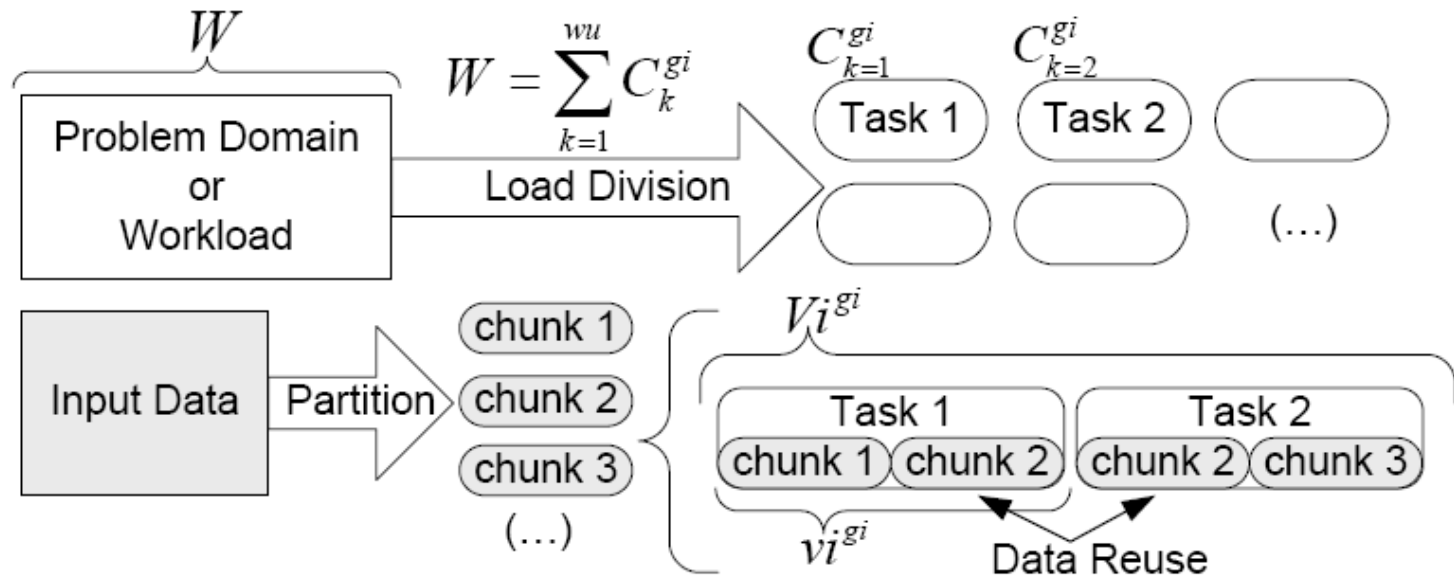
5. Conclusions



Assumptions

- ▶ Tuning on frameworks provides a balance between performance improvement and generalization.
- ▶ Master–Worker programming helps application deal with system heterogeneity.
- ▶ The significant parameters for Master–Worker tuning is Number of Workers and Grain Size (Compute/Computation ration) selection.
- ▶ Application should support change of number of workers and grain size selection at runtime.

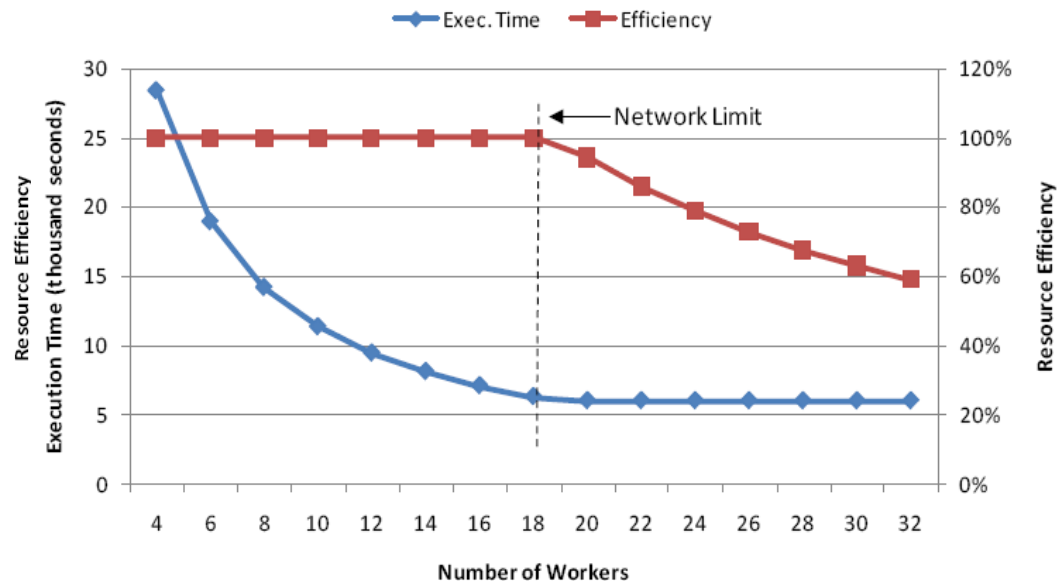
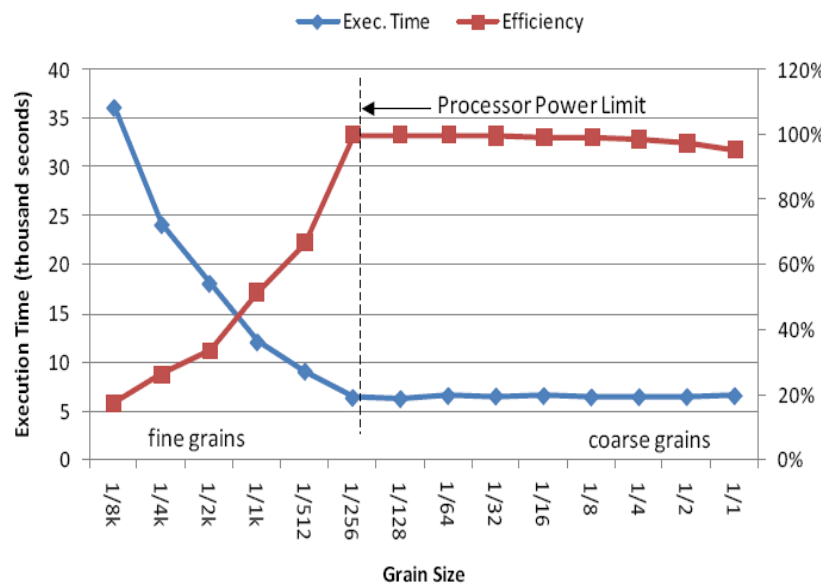
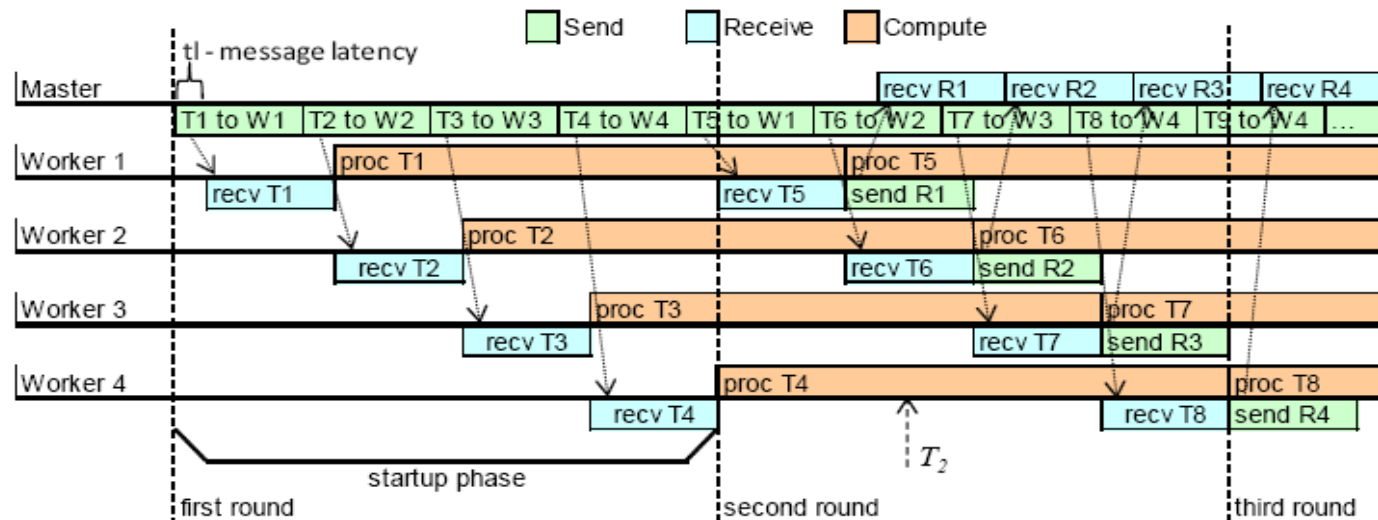
Changing Grain Size (Compute/Communication Ratio)



- ▶ Smaller Tasks → Higher communication volume
- ▶ Coarser Tasks → Lower communication volume

**Grain Size Selection
changes Compute/Communication Ratio**

Master/Worker Applications



Performance Indexes

▶ Scenario and Problem

- Given a set of resources assigned and mapped by a meta-scheduler/resource broker, how to reduce application execution time and increase resource usage efficiency?

▶ Objectives

- Lower total execution. ← Grain size selection.
- Increase resource usage efficiency. ← Number of Workers.

Performance Models

- ▶ The minimal execution time is achieved when application have the maximum number of faster workers in busy state.
- Grain Size selection helps application to use the assigned resource set.

$$N_{opt}^{gi} = \frac{Tc^{gi}}{\max(vi^{gi}, vo^{gi}) * \lambda}$$

Direct impact on
resource usage
efficiency

$$\left\{ \begin{array}{l} T_{startup}^{gi} = tl + (vi^{gi} * \lambda) * \left(\frac{N_w + 1}{2} \right) \\ T_{finalization}^{gi} = tl + (vo^{gi} * \lambda) * \left(\frac{N_w + 1}{2} \right) \end{array} \right.$$

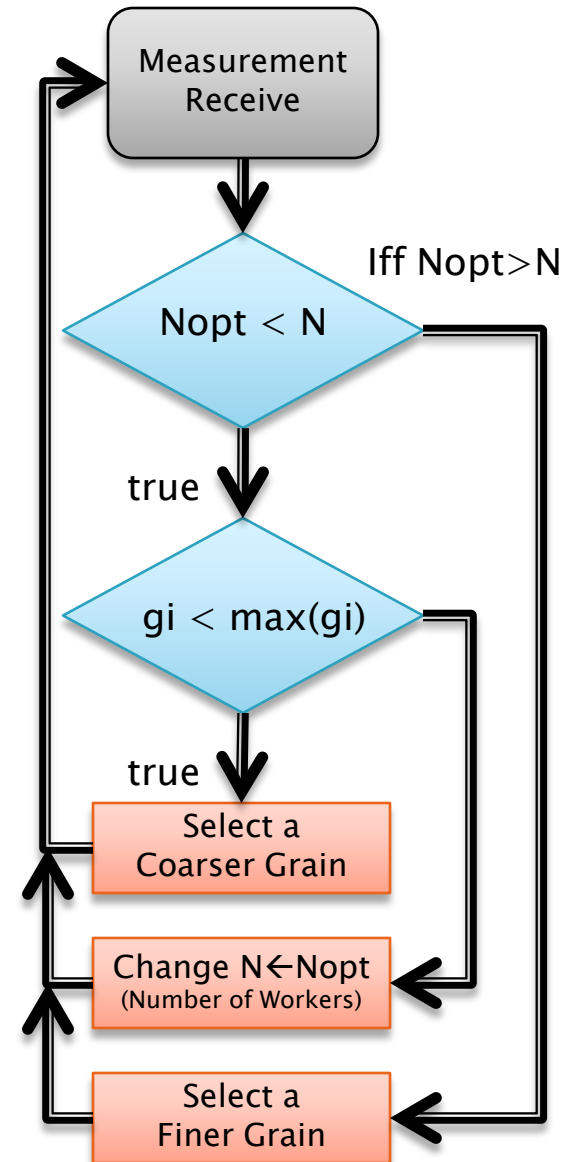
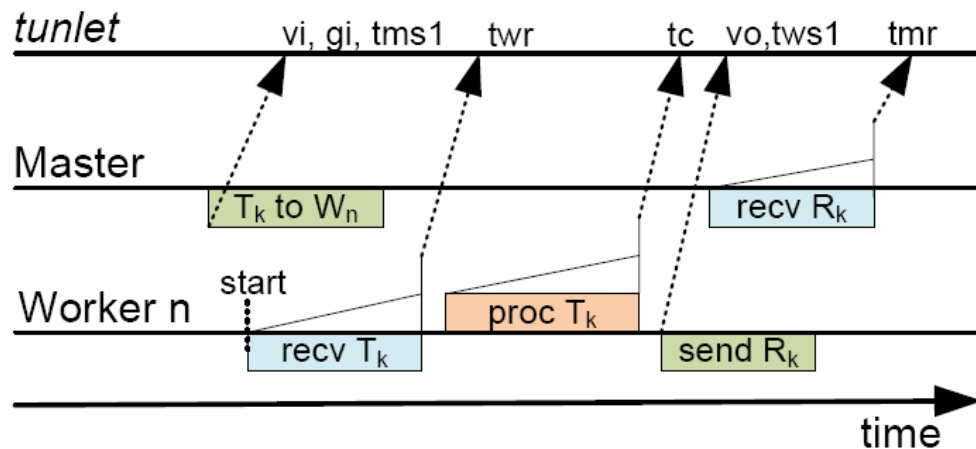
Content

1. Motivation
2. Performance Model for Dynamic Tuning
3. Tuning Heuristics
4. Experimentation Results
5. Conclusions



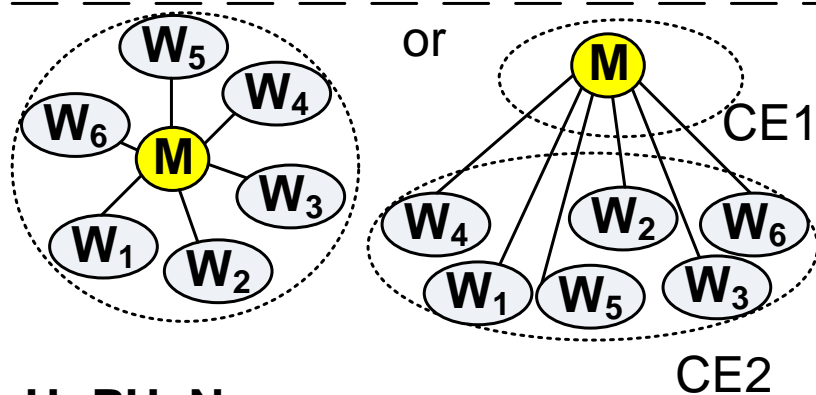
Tuning Process

- ▶ Runtime compute and communication of each task per worker are continuously measured.

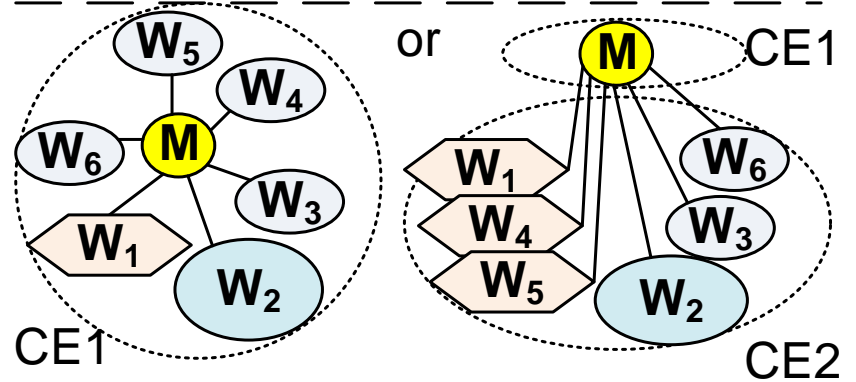


Processor and Network Heterogeneity

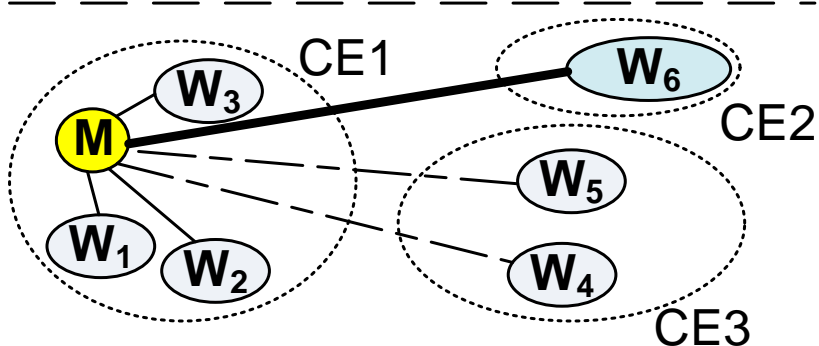
HoPHoN



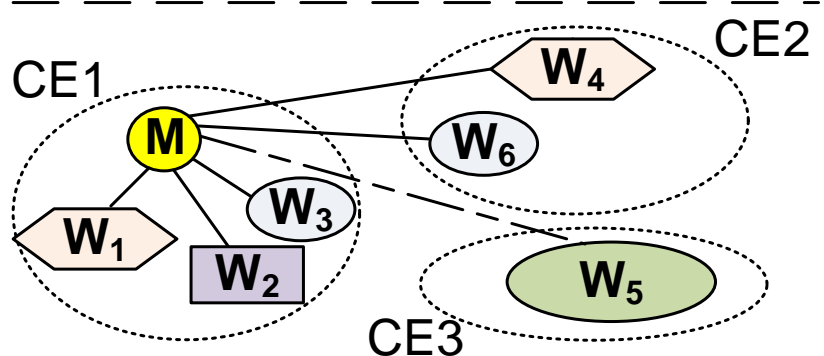
HePHoN



HoPHeN



HePHeN



Heterogeneous Scenarios

- ▶ Task distribution gives higher priority to faster processors.
- ▶ Remote Compute Hosts (CH) are grouped by Compute Element (CE) which receives tasks in parallel.
 - Communication managers isolates master to remote worker communication blocking.
- ▶ The number of workers in each CE is independently tuned (greedy approach).

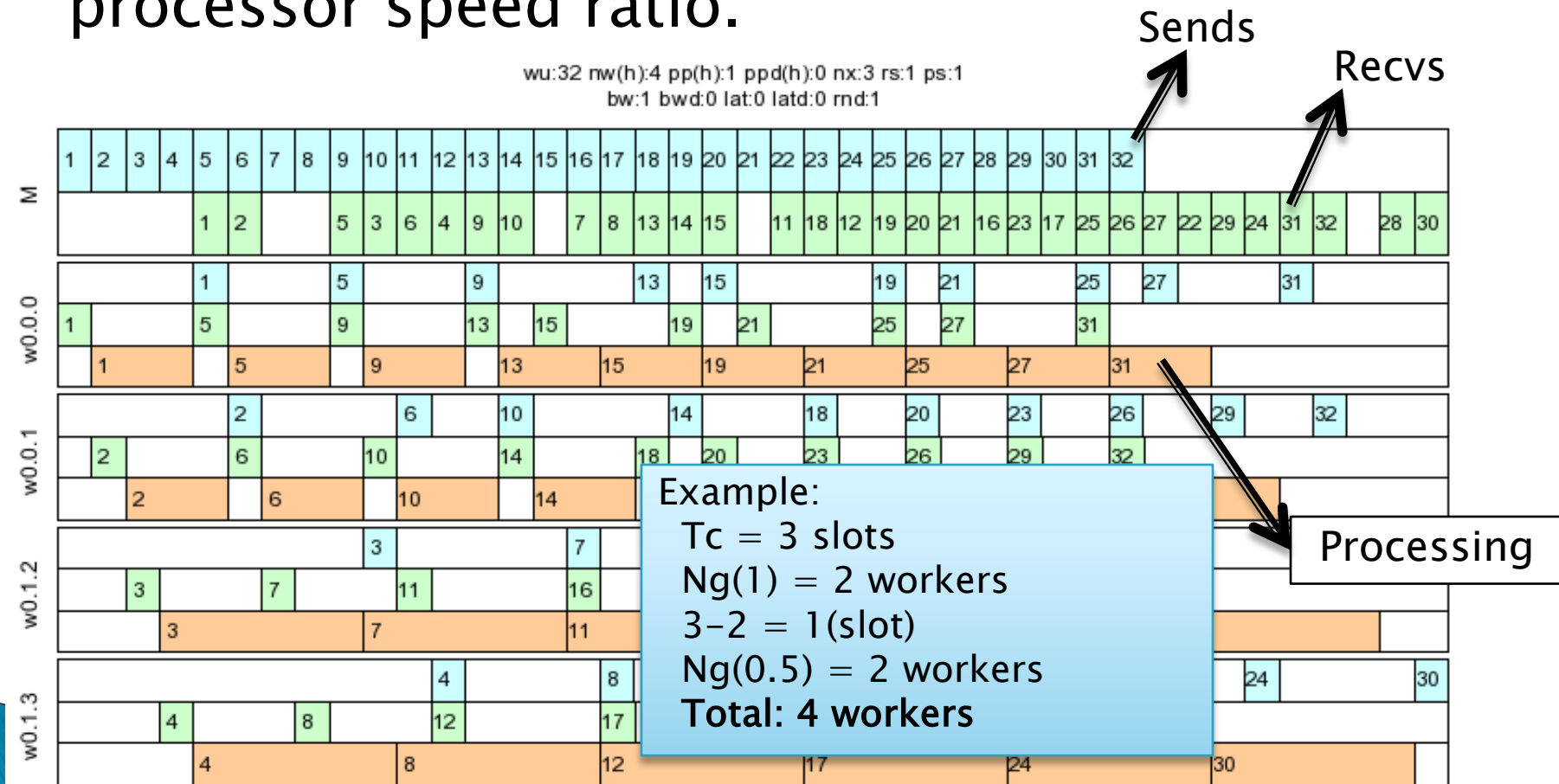
Tune Heuristic



```
let slot = 0 and nw = 0
for each Wg(k) from k=0 to Wgn-1
  let Nopi = calc(Nopt for Wg(k)) # acquire measurements
  if Nopi <= Wg(k) then
    if gi>0 then
      let gi = gi + 1 #tune grain size
    else
      let Nopt = Nopi # tune num.workers
    end if
  exit for
else
  let Tc(k) = (Tc(0)-slot) * (Tc(k)/Tc(0))
  let Nopi = calc(Nop in Tc(k))
  if Nopi < Wg(k) then
    let Nopt = Nopi # tune num.workers
    exit for
  else
    let Nw = Nw + Wg(k)
    let slot = slot + Wg(k) * (Tc(k)/Nopi))
  end if
end if
end for
```

Number of Workers with Heterogeneous Processors

- Allocation of communication time slot using processor speed ratio.



Content

1. Motivation
2. Performance Model for Dynamic Tuning
3. Tuning Heuristics
4. Experimentation Results
5. Conclusions



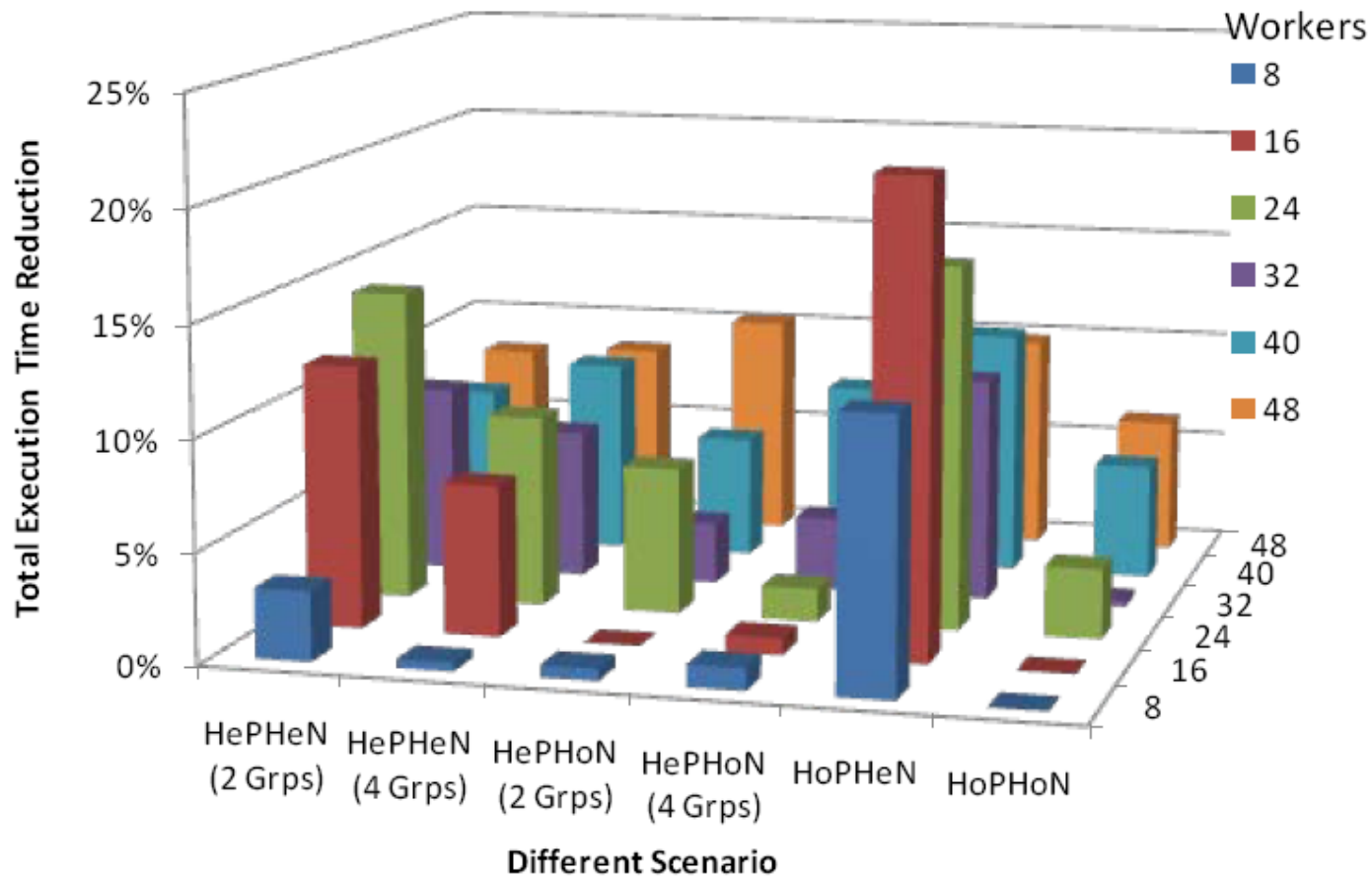
Experimentation Methodology

- ▶ Application execution may have different measurements in different executions.
- ▶ Cases coverage:
 - If grain size and number of worker tuning do not raise execution time in all different execution scenarios and system characteristics change, the new scenario found should also have benefits from the tuning.
- ▶ Event driven simulation using real task execution times and network bandwidth.

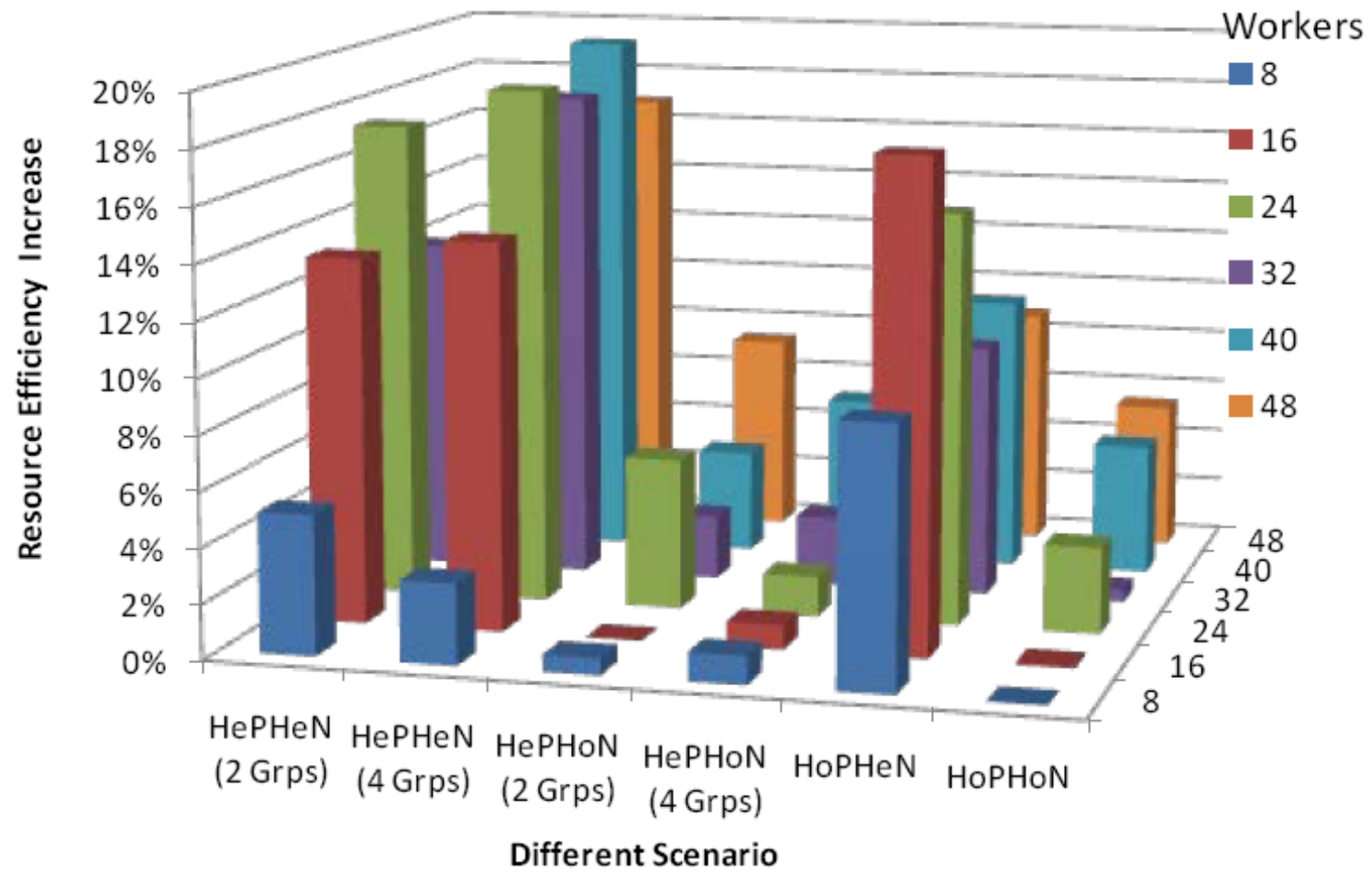
Workload

- ▶ Matrix Multiplication Application with block row and column division parallelization strategy with dynamic load balance (faster processors first).
- ▶ Simulation parameters
 - Task Execution time → ATLAS sgemm
 - Network Bandwidth
 - Local Workers → MPI
 - Remote Workers → iperf tool
- ▶ Cases with processor and network heterogeneity and grain sizes from compute to communication bound scenarios.

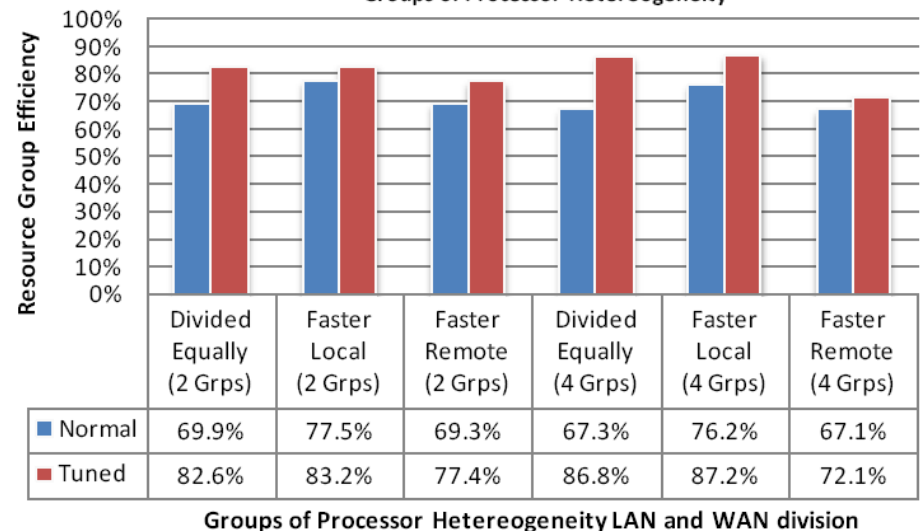
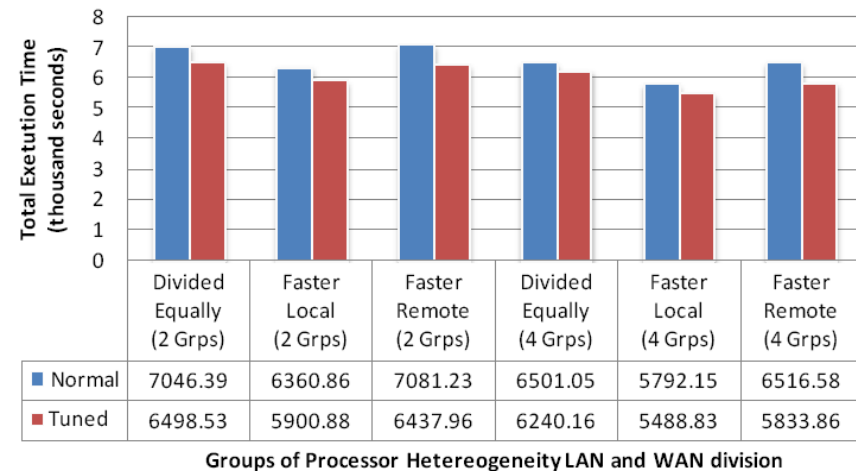
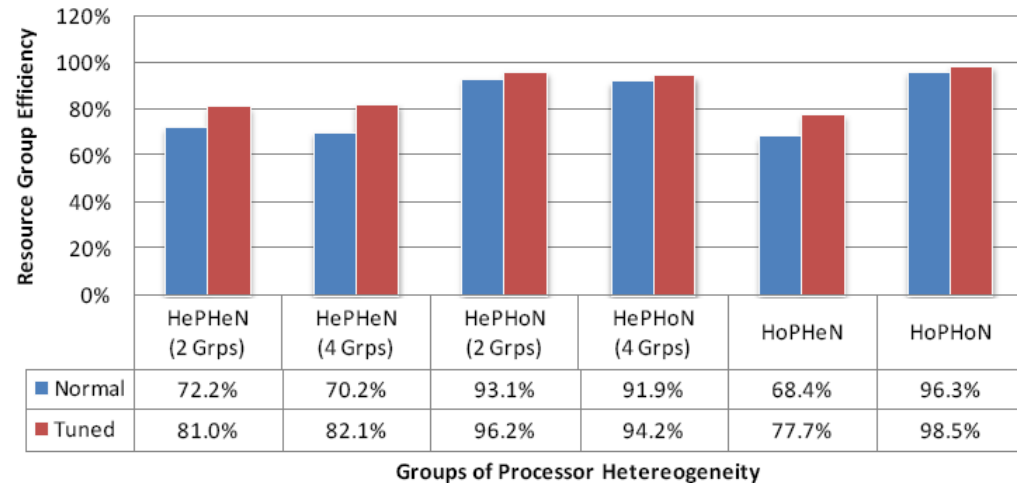
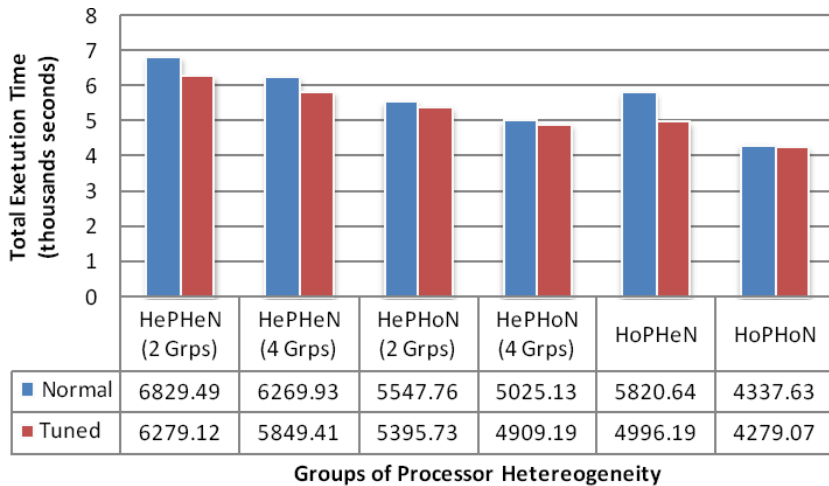
Total Execution Time Analysis



Resource Usage Efficiency Analysis



Heterogeneity Analysis



Content

1. Motivation
2. Performance Model for Dynamic Tuning
3. Tuning Heuristics
4. Experimentation Results
5. Conclusions



Conclusions

- ▶ We proposed a heuristic to dynamically change application grain size and number of workers.
 - Deal with network and processor heterogeneity.
 - Increase application performance indexes.
- ▶ The heuristic was evaluated over a Master/Worker application in different processor and network heterogeneous scenarios.
 - Computational Grids composed by commodity clusters scattered over Internet .
- ▶ The best gains where obtained on scenarios with higher level of heterogeneity.

Open Lines

- ▶ Apply our heuristic on different Grid master/worker applications.
- ▶ Use same heuristic to tune pipeline stage replication in pipeline execution over Grids.
- ▶ Grain Size selection on SPMD applications.
- ▶ Hierarchical dynamic tuning of grain size and number of workers.

Thank You!

Questions??