# Scheduling Divisible Tasks with Message Passing Interface
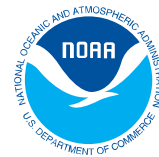
By: Jaime Ballesteros, PhD Student

Advisor:

Prof. Jaime Seguel

Parallel and Distributed Computing Laboratory

University of Puerto Rico at Mayaguez (UPRM)

May 2007

# Problem formulation

How to effect task scheduling in a distributed system, specifically in numerical simulations with an MPI-flavor, in order determine the most effective orchestration of communications and computations that will give the optimal throughput of the system.

**Justification:**

Numerical simulations require optimal throughput in order to return accurate results in a reasonable time. Sometimes using a distributed system to compute simulations will give better results, but orchestration of communications and computations has to obey a predefined scheduling policy, that, in some cases, is unlikely to reach an optimal throughput in a expected time.

| ID | Processor Name | Duration | Time Units | | | | | | | | | | | | | |
|----|----------------|----------|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| | | | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 |
| 1 | Processor 1 | 2d | | ███ | | | | | | | | | | | | |
| 2 | Processor 2 | 1d | | | | | | | | | | | | ███ | | |
| 3 | Processor 3 | 4d | | | | ████████ | | | | | | | | | | |
| 4 | Processor 4 | 3d 4h | | | | | | ████ | | | | | | | | |
| 5 | Processor 5 | 4d | ████ | | | | | | | | | | | | | |
| 6 | Processor 6 | 2d 4h | ██ | | | | | | | | | | | | | |
| 7 | Processor 7 | 3d | | | ██████ | | | | | | | | | | | |
| 8 | Processor 8 | 4d | | | | | | | ██████ | | | | | | | |
| 9 | Processor 9 | 1d | | | | ██ | | | | | | | | | | |
| 10 | Processor 10 | 4d 4h | | | | | | ██████ | | | | | | | | |

# Methodology (Solution Approach)

**Steady-state scheduling:**
• Provides the asymptotically optimal throughput scheduler in master-slave applications (i.e. divisible applications).
• Solved in polynomial time with linear programming

Maximize $n_{\text{task}}(G) = \sum_{i=1}^{p} \dfrac{\alpha_i}{w_i}$

Subject to

$$\forall i, \qquad\qquad 0 \leq \alpha_i \leq 1$$

$$\forall i, \forall j \in n(i), \qquad 0 \leq s_{ij} \leq 1$$

$$\forall i, \forall j \in n(i), \qquad 0 \leq r_{ij} \leq 1$$

$$\forall e_{ij} \in E, \qquad\quad s_{ij} = r_{ij}$$

$$\forall i, \qquad\qquad \sum_{j \in n(i)} s_{ij} \leq 1$$

$$\forall i, \qquad\qquad \sum_{j \in n(i)} r_{ij} \leq 1$$

$$\forall e_{ij} \in E, \qquad\quad s_{ij} + r_{ij} \leq 1$$

$$\forall i \neq m, \qquad \sum_{j \in n(i)} \dfrac{r_{ij}}{c_{ij}} = \dfrac{\alpha_i}{w_i} + \sum_{j \in n(i)} \dfrac{s_{ij}}{c_{ij}}$$
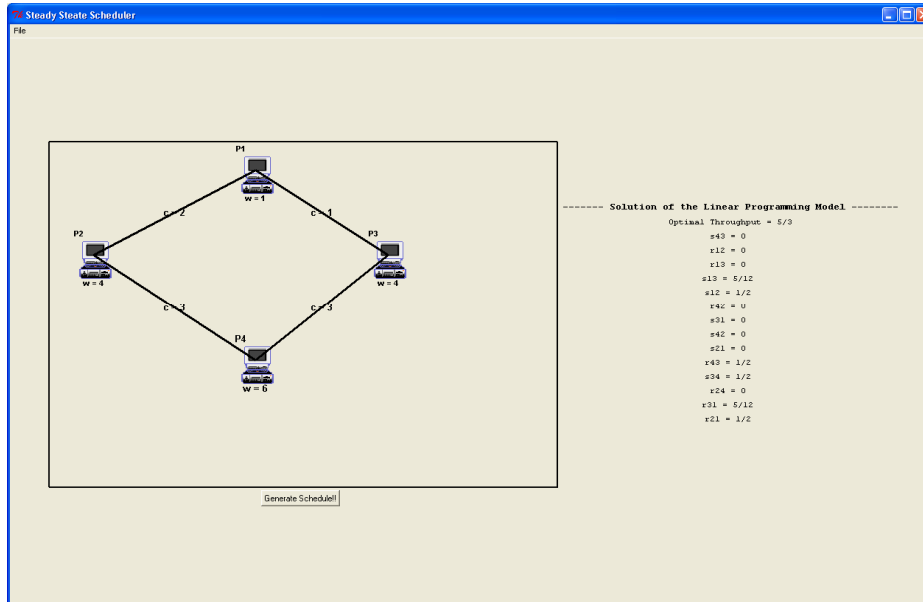
$$\forall i \in n(m), \qquad r_{mj} = 0$$

**Elements of divisible task complexity theory:**
• Atomic tasks, unit tasks; volume and communications
• Provides the asymptotically shortest schedule for processes with single communication phase. All process end processing at the same point in time.
• Solved in polynomial time

**The system will receive the user code, identify the atomic tasks and communication graph, applied the theoretical framework and emit code, pretty much in the spirit of the FFTW or some parallel data base search algorithms.**

WALSAIP

P D C Group
Parallel &
Distributed
Computing  6

# Applications Tools



**In our first study of Steady State scheduling theory we implement Steady State Scheduler V1.0 in Python and it:**

- Allows to change communication and execution times.
- Uses Glpk® to solve the Master-Slave linear programming problem.
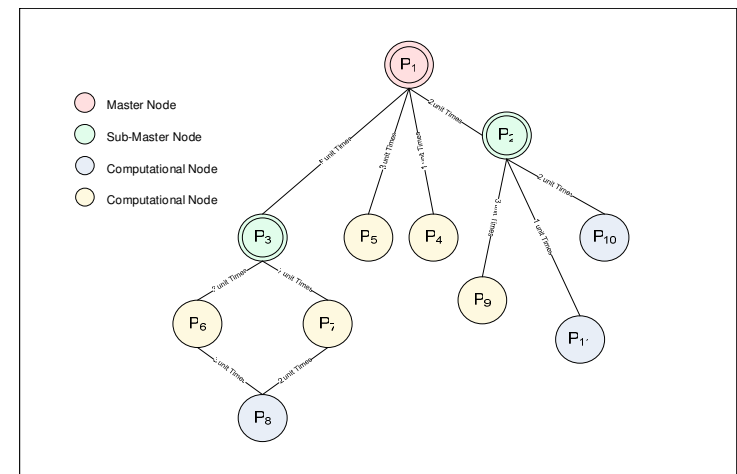- Constructs a theoretical schedule.

**Open MPI and GCC will be our primary tools:**

Process and communication scheduling will be reflected in the source code and "controlled" with a feed-back control mechanism based on the communication flow.

# Research Results

1. By using the demo, we could identify some observables in the system and the subsequent behavior in order to make an in-depth study of steady state scheduling mechanism.
2. We have extended (or perhaps) unified the theory of load divisible and steady state scheduling for application tasks that can be mapped as starts or trees.
3. We first modified the divisible load scheduler to make it periodic, saving thus a significant amount of start-up overhead. Then, we applied the same technique to the steady state scheduler to get a hybrid method that is provable superior to both of its ancestors.
4. We also developed a communication centric formulation of the new scheduler. Such formulation allows for the absorption of transients (decline in the processor or network speed).