


# Scheduling Message Passing

(what Meta-Computing systems need to do)



CS-594  
Dr Graham E. Fagg  
Spring 2001

## Overview

- Covers details you may have missed in David Walkers MPI / Message Passing class.
- Assumes that you know and understand MPI and PVM.

CS-594 Scheduling Message Passing applications

## Closely coupled vs Cluster Computing

- Bottom line
  - MPI is better at message passing than PVM
  - More complex
  - Less flexible at anything else
    - I.e. it's a message passing system not a distributed environment

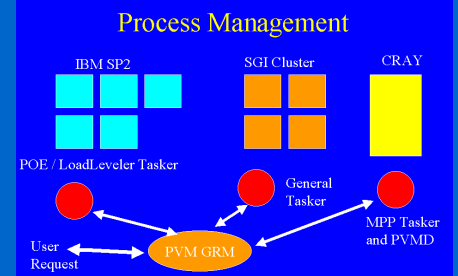
CS-594 Scheduling Message Passing applications

## Scheduling

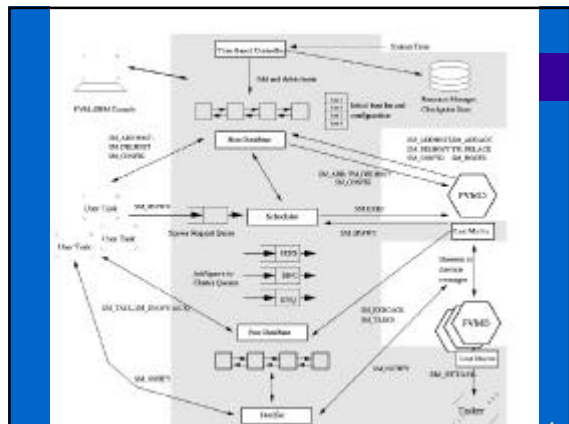
- Not built into MPI as it has no process control
  - But maybe an option under MPIRUN
  - Yep PVM has it all
    - user controllable pvm\_spawn()
    - pvm\_rm interface
      - also
        - pvm\_tasker interface
        - pvm\_hoster interface

CS-594 Scheduling Message Passing applications

## Scheduling Process Management



CS-594 Scheduling Message Passing applications



## Scheduling

- Two types to worry about
  - At spawn time
    - static allocation based on the environment
  - At run time
    - I.e. migration of tasks
      - system level migration
        - Special support needed (Condor)
      - User level
        - check points / restarts
  - Change work load allocated (bag of tasks)

CS-594 Scheduling Message Passing applications

## Task allocation in PVM I.e. pvm\_spawn()

- Before improving on it, had to figure out how it worked as it wasn't random but round-robin
- Aimed at using spare capacity
  - what spare capacity??

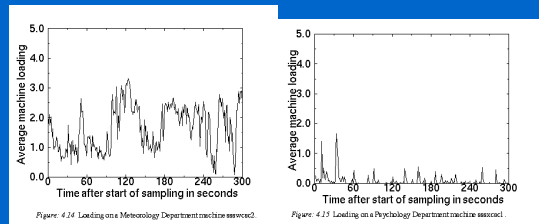
CS-594 Scheduling Message Passing applications

## What is spare (what is even machine load?)

- Condor people claimed 10% utilisation for their systems
  - At Reading was more like 40-60% all the time.
- Load
  - machine average is not a good metric but without more specific help from the kernel it would have to do.
- Defined user classes and loading based on observations on the RDG system over a year...

CS-594 Scheduling Message Passing applications

## Typical loading



CS-594 Scheduling Message Passing applications

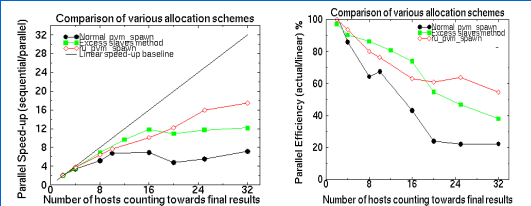
## Better Spawn

- Added checks for load before starting.
  - Based on two methods, central RPC and distributed monitor daemons
  - Checked effects of this system on startup performance, and accuracy of placement.

Scenario	1 (on P1)	70	average	Stago
normal	10.000	10.000	10.000	10.000
pvm_spawn	1.346	4.4	1.1	0.045
central_spawn	1.141	31.5	14.4	0.251
distributed_spawn	1.648	17.4	3.8	0.294

• Researcher has to test the use of static-allocated loads. He has to test the parallelism and the accuracy of placement with the use of the file and the distributed system to avoid the use of the file.

## Application Performance



CS-594 Scheduling Message Passing applications

