# Coordinated Coscheduling in Time-Sharing clusters through a Generic Framework

**Saurabh Agarwal**

**IBM Research (India Research Labs)**

**IIT New Delhi – 110016, India**

**Co-authors:**

Gyu Sang Choi  : PhD Student, CSE Dept. , Penn State University.

Dr. Chita Das   : Prof. (& my Advisor), CSE Dept. , Penn State University.

Dr. Andy B. Yoo  : LLNL, Livermore, CA

Dr. Shailabh Nagar: IBM T.J. Watson research center, NY

PENNSTATE

# Clusters Today

Scientific Applications

Web / Wireless Application Servers

# Typical cluster based Application Characteristics

Parallel and or Distributed.

Possibly Communication Intensive (High/Medium/Low).
- Very common and is our research focus.

Requirement to progress ?
- Timely Inter-Node communication.

Implication ?
- Parallel Jobs must be co-scheduled !

# Outline Today

- What is the **co-scheduling** problem ?

- How has been solved **earlier** ?
    - Batch, Gang, DCS, ICS (SB), PB

- Why prior solutions **not enough** ?
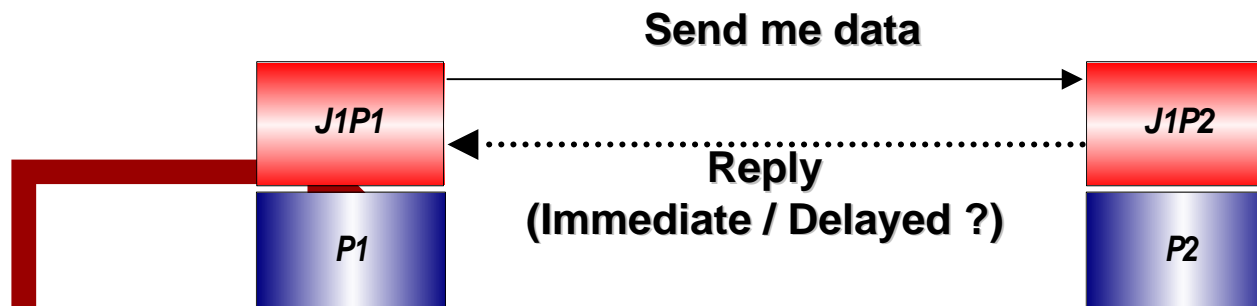
- What do I **propose** ?

- What are the **results** ?
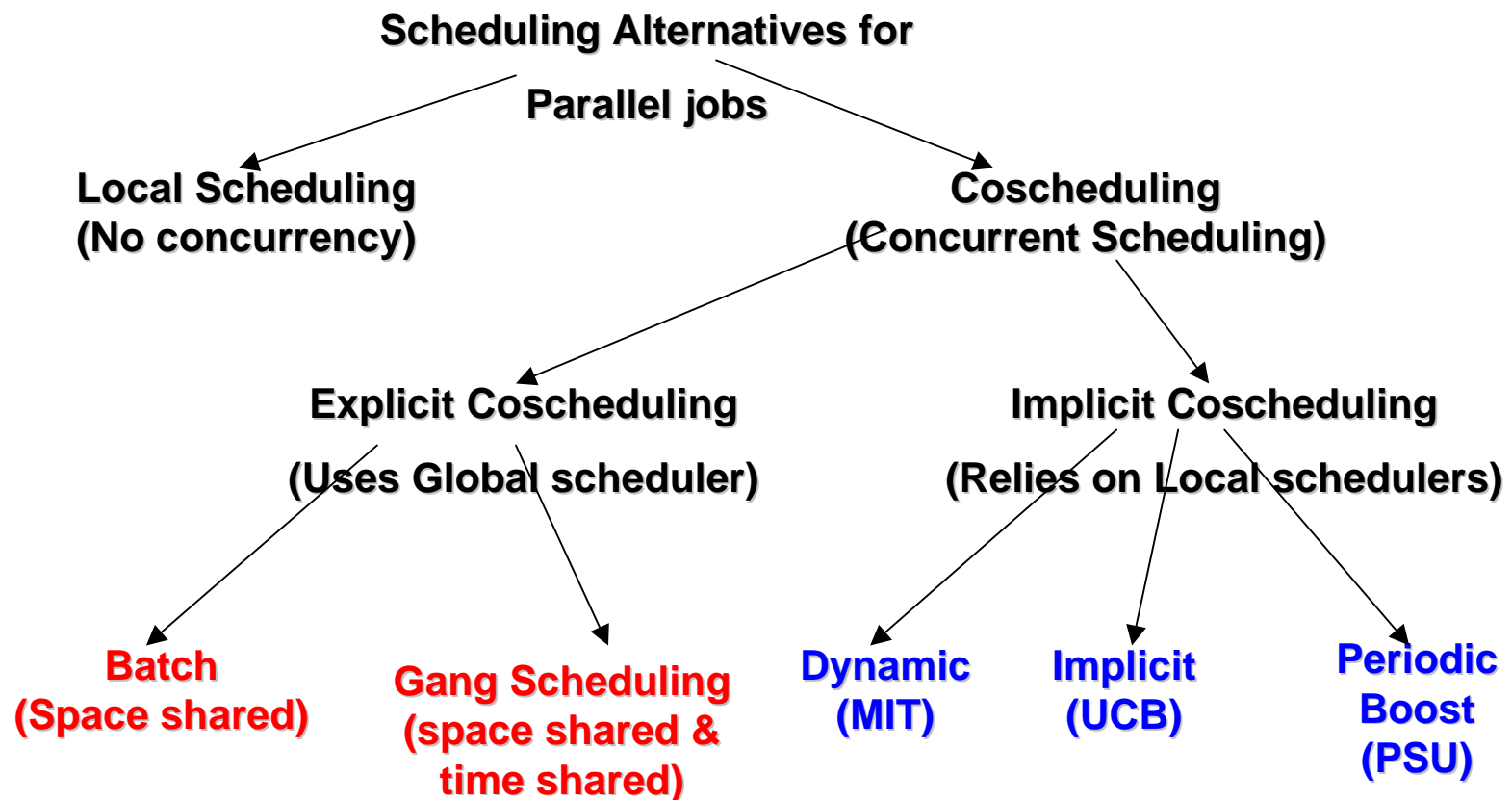
- What can we **conclude** ?

- Is there a **future** for this ?

# What is CoScheduling ?

CoScheduling* :  Concurrently Schedule processes of a parallel job on individual nodes of a time-sharing cluster.
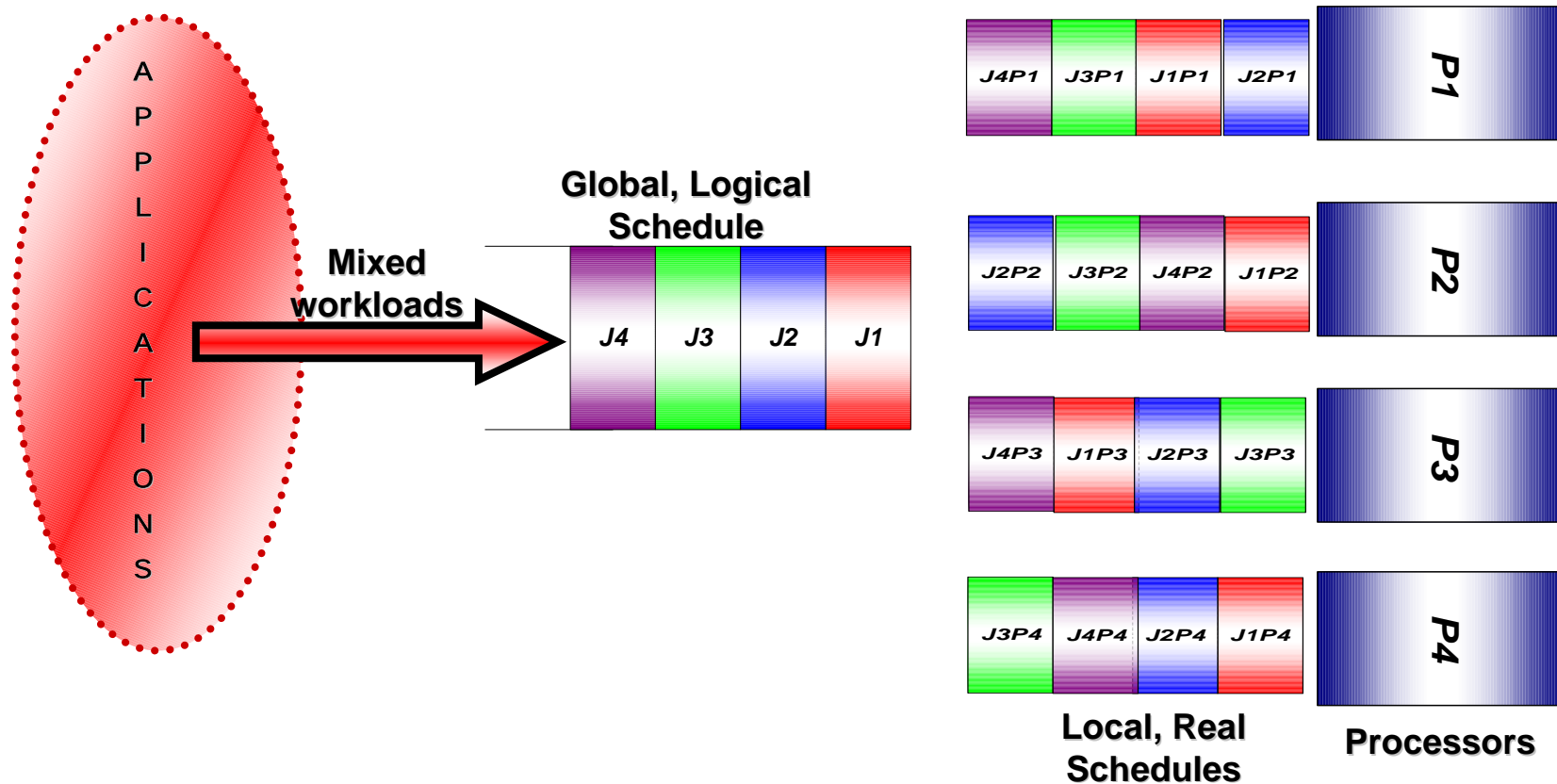
**Send me data**

| J1P1 | ➤ | J1P2 |

**Reply (Immediate / Delayed ?)**

| P1 | | P2 |

Why Needed ? :  **Performance**

PENNSTATE

# Parallel job Scheduling Techniques (Hierarchy)

**Scheduling Alternatives for**

**Parallel jobs**

**Local Scheduling**
**(No concurrency)**

**Coscheduling**
**(Concurrent Scheduling)**

**Explicit Coscheduling**

**(Uses Global scheduler)**

**Implicit Coscheduling**

**(Relies on Local schedulers)**

**Batch**
**(Space shared)**

**Gang Scheduling**
**(space shared &**
**time shared)**

**Dynamic**
**(MIT)**

**Implicit**
**(UCB)**

**Periodic**
**Boost**
**(PSU)**

# Un-coordinated (Local) scheduling



APPLICATIONS

Mixed workloads

**Global, Logical Schedule**

| J4 | J3 | J2 | J1 |

| J4P1 | J3P1 | J1P1 | J2P1 | **P1** |
| J2P2 | J3P2 | J4P2 | J1P2 | **P2** |
| J4P3 | J1P3 | J2P3 | J3P3 | **P3** |
| J3P4 | J4P4 | J2P4 | J1P4 | **P4** |

**Local, Real Schedules**

**Processors**

# Explicit coscheduling: Batch Scheduling

- 1 job runs on all CPUs until completion
- Other jobs wait in scheduler queue for their turn
- Problems
  - **Low utilization**
  - **Low response times**
- Examples :
  - IBM-SP2 (Uses Load Leveler)
  - Intel Paragon (Uses Network Queuing System (NQS) )
  - Many research COTS clusters use PBS

# Explicit Coscheduling: Gang Scheduling

**A P P L I C A T I O N S**

**Mixed workloads**

**Global, Logical Schedule**

| J4 | J3 | J2 | J1 |

Examples :

IBM-ACSI White

IBM ASCI Blue-Pacific

| J4P1 | J3P1 | J2P1 | J1P1 | **P1** |

| J4P2 | J3P2 | J2P2 | J1P2 | **P2** |

| J4P3 | J3P3 | J2P3 | J1P3 | **P3** |

| J4P4 | J3P4 | J2P4 | J1P4 | **P4** |

**Local, Real Schedules**    **Processors**

# Why do Gang Scheduling on Clusters ?

Need a single scheduler controlling jobs on ALL nodes

- Impractical to implement as :-

  - Frequent Synchronization (Order milli-sec)
    - Costly in **NOW** because of higher wire latencies (Remember ?).

  - Requires large 'time quantum' (of order seconds).
    - Reduces system responsiveness.

Not scalable for NOW (For same reason as above).

**PENN**STATE

# What is the solution ?

**Gang Scheduling**    **Implicit Co-scheduling**

| P4 | J4P1 | J3P1 | J2P1 | J1P1 |
| P3 | J4P2 | J3P2 | J2P2 | J1P2 |
| P2 | J4P3 | J3P3 | J2P3 | J1P3 |
| P1 | J4P4 | J3P4 | J2P4 | J1P4 |

J4P1 J3P1 J2P1 J1P1
J4P2 J3P2 J2P2 J1P2
J4P3 J3P3 J2P3 J1P3
J4P4 J3P4 J2P4 J1P4

**Co-scheduled only for part of the time slot.**

**Time**

PENNSTATE

# DCS : Dynamic Coscheduling

**Real Applications, || Workloads**

**Middleware (MPI / VIA)**

**OS Kernel**

**Scheduler**

**4**

**NIC Driver**

**3**

**2**

**DMA**

**NCP**

**1**

**Intelligent NIC**

**NIC Processor**

*"Dynamic Adjustment"* of process priorities by using "*messages*" (from communicating jobs) as hints.

**NIC Control Program (NCP) senses incoming message.**

**NCP raises interrupt(s) to NIC Driver, *if required*.**

**NIC Driver *finds* & *places* relevant process into the highest-priority queue.**

**Scheduler schedules that process.**

PENNSTATE

# in Block

# PB : Periodic Boost

**Real Applications, || Workloads**

**Middleware (MPI / VIA)**

**OS Kernel**

**Scheduler**

**1**

**3**

**2**

**NIC Driver**

**DMA**

**Intelligent NIC**

**NCP**

**NIC Processor**

*"Periodically"* check end-points in host memory, "**Boost**" the process with pending messages.

**Driver p   s VI end-points in user-pin   d memory, every 10 ms.**

**Picks 1   ending-message VI , finds co   esponding process, boosts i   riority.**

**Scheduler schedules that process next.**

PENNSTATE

# 'Implicit' Coscheduling: Various Issues

**DC...**

- T... ...terrupts (< 1ms )
- R... ... VI, instead of per process

**SB...**

- N... ...ost...
- N... ...de...de optimizatio...
- V... ...t ... For Tightly coupled (Split-C) environment.

**PB...**

- Un... ...rent form.
- De... ...uracy : Polling done after DMAs
- Lo... ...ing in fair version. (Too many polls)
- Reason : Per VI, instead of per process.

Why prior solutions **not enough** ?

# No Commercial Implementation. WHY ?

1. Lack of exhaustive experimentation on multiple platforms.
   - Not easy to code custom-solutions for each platform.
   - No general standard approach available.

2. No scheme best for all types of workloads.
   - None provides extensibility, generality, adaptability ?
   - Support for QoS not addressed at all.

3. No real incentives (results) demonstrated yet
   - Coscheduling against batch scheduling ?
   - Presence of other sequential workloads (CPU, I/O) ?
   - High Multi-programming degree ?

# Addressing global issue 1 :- Prior Design

**Real Applications, || Workloads**

**Middleware (MI**

**Communication Libra**

**OS Kernel**

**Scheduler**

**NIC**

**Intelligent NIC**

*Parts of code to change :*

*MPI Library (SB)*

*VIA Library (SB, PB)*

*Device driver (DCS, SB, PB)*

*Firmware (DCS, SB, PB)*

**Disadvantages**

**Flexibility :** Tight driver/firmware coupling

**Generality :** Different implementations for each scheme.

**Modularity :** No re-use across platforms.

**Portability :** No standard interface

**Integrity :** Local scheduler isolated.

# Proposed Design

**Real Applications, || Workloads**

**Middleware (MPI)**

**Communication Library (VIA)**

**coschedLib**

**Cosched Module**

**Alt. sched Patch**

**Scheduler**

**NIC Driver**

**OS Kernel**

**Intelligent NIC**

**CP**

**Big Picture :**
- **Identify and isolate "*policies*" and "*mechanisms*".**
- **Re-use "*policies*"**
- **Re-implement "*mechanisms*"**

**Advantages**

**Flexibility :** Coscheduling independent of device driver / firmware.

**Generality :** Same module for all schemes.

**Modularity :** Re-use across platforms.

**Portability :** Standard interfaces defined.

**Integrity :** Local scheduler involved.

PENNSTATE

# Addressing global issue 2 :-
# Coordinated Coscheduling

# Addressing global issue 3 :- Workload / Environment

16 node Myrinet connected cluster, 1GB RAM .

Linux 2.4, MPI over VIA,  Berkeley-VIA

Lanai-9 Myrinet NI cards, 8MB on-chip memory.

| Workload | Applications | Communication Intensity |
|---|---|---|
| $Wl1$ | (EP, EP, EP, EP, EP, EP) | lo:lo:lo:lo:lo:lo |
| $Wl2$ | (EP, EP, EP, MG, MG, MG) | lo:lo:lo:hi:hi:hi |
| $Wl3$ | (MG, MG, MG, MG, MG, MG) | hi:hi:hi:hi:hi:hi |
| $Wl4$ | (EP, EP, LU, LU, MG, MG) | lo:lo:me:me:hi:hi |
| $Wl5$ | (LU, LU, LU, LU, LU, LU) | me:me:me:me:me:me |

(a) Parallel Workload Composition (MPL6

| Category | Workload Mix |
|---|---|
| Parallel Only | $wl1, wl2, wl3, wl4, wl5$ |
| Parallel + CPU | $(wl1...wl5) + 1\ sb$ |
|  | $(1,2,4,6)sb + 2$ MGs |
| Parallel + IO | $(wl1...wl5) + 1\ iobench$ |

(b) Executed Combinations ($sb : sched\_bench$)

Table   : **Workload mixes used in this study.**

# Performance : Execution Time

**MPL = 6**

CC > SB > PB > DCS > Local

# Where does time go ?

Total CPU Time = Useful work time + Non-useful spin time

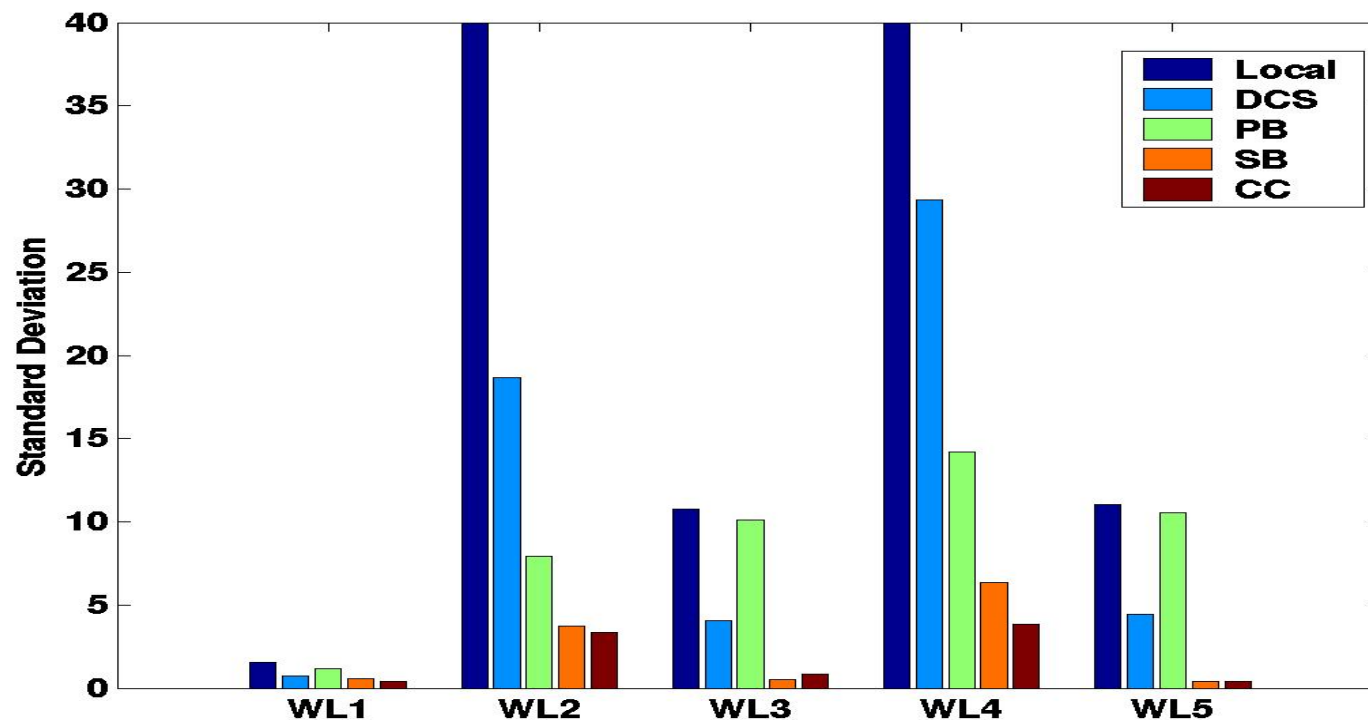Wall Clock Time = Total CPU time + Block time + others...



**Workload –** *wl3*

**Workload –** *wl4*

# Tolerance : Standard Deviation

*Wl2, Wl4* : Mixed                    *Wl1, Wl3, Wl5* : Uniform
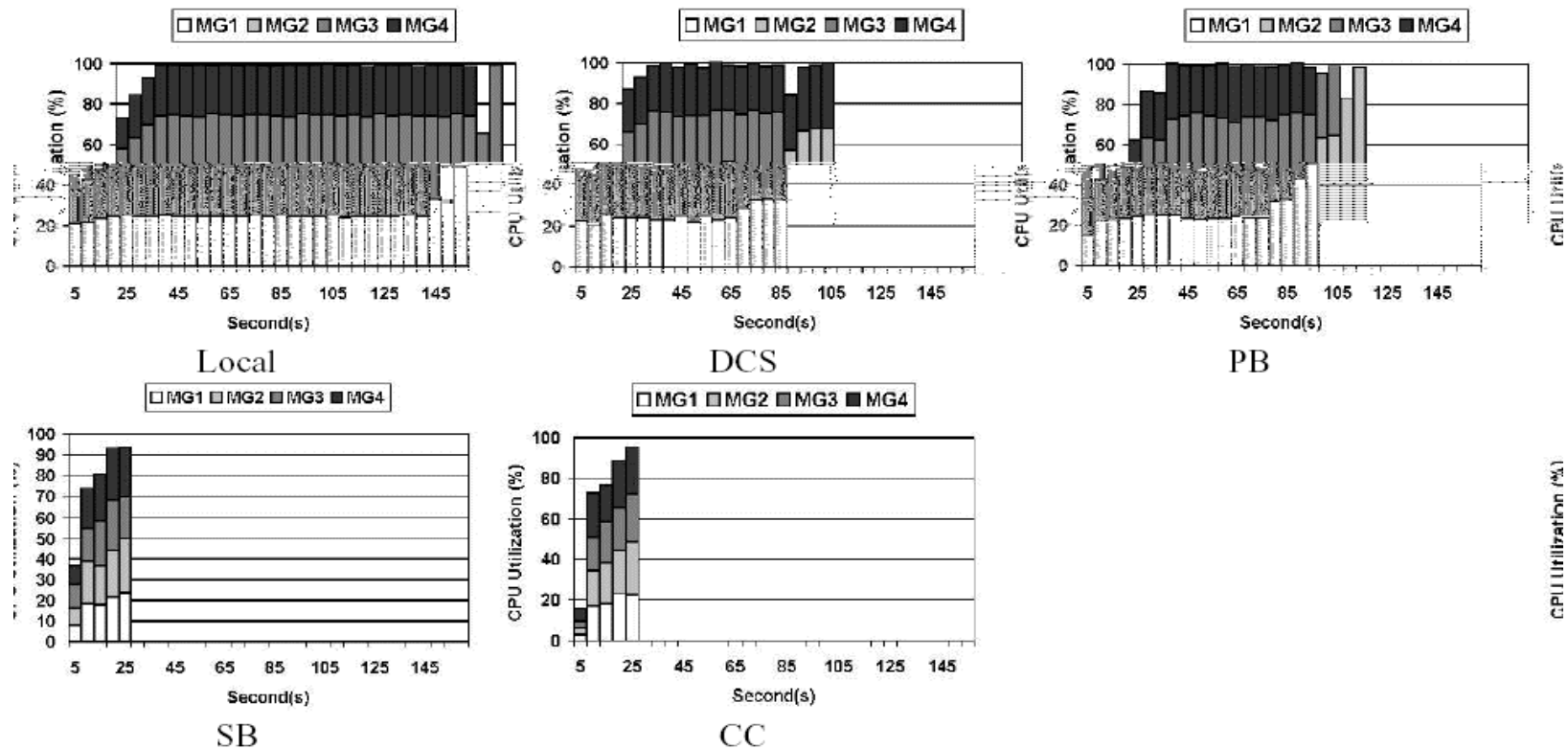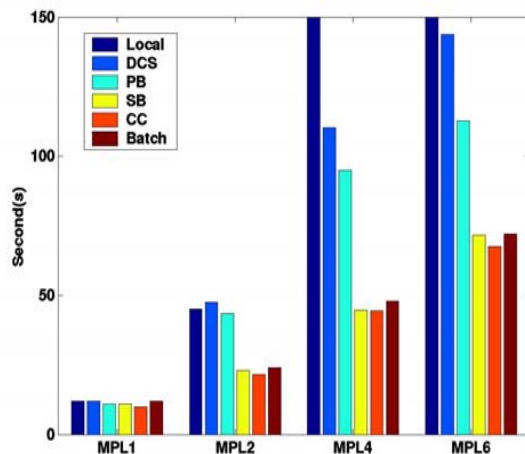
# Fairness : CPU Utilization



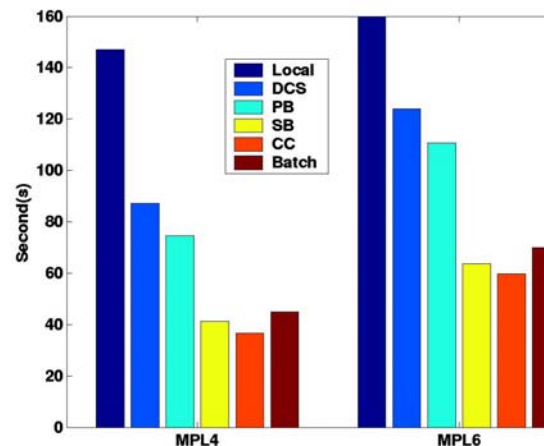Figure : CPU Utilization results for wl3 (MPL=4)

# Scalability (vertical)
# Effect of MPL
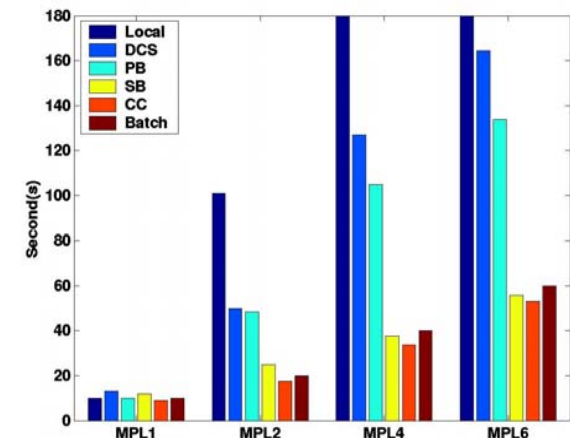
Rate of increase higher in Local, DCS, PB than SB, CC.

Glimpse comparison to batch scheduling.



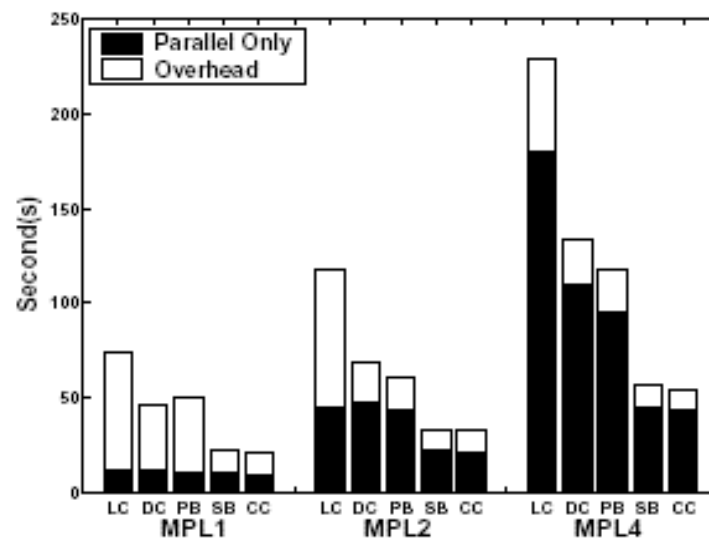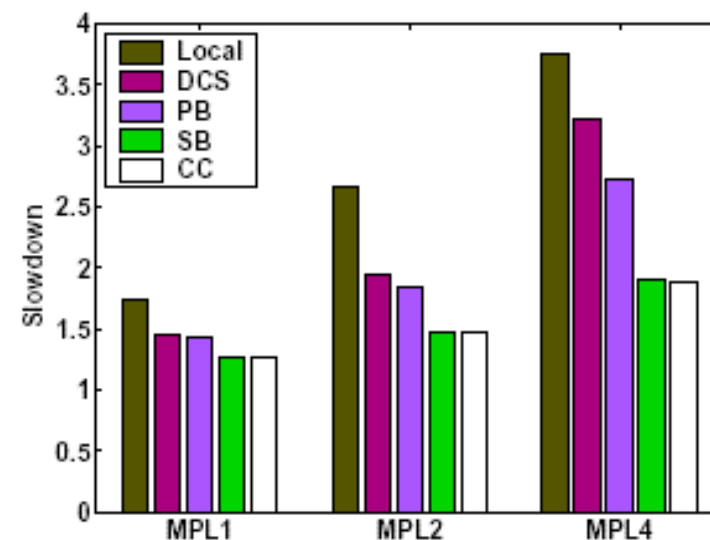**Workload –** *wl3*          **Workload –** *wl4*          **Workload –** *wl5*

# Mixing CPU intensive jobs (W*I3*)

CC & SB tolerate load better (low *overhead* in (a) ).
CC and SB exploit idle cycles well (low *slowdown* in (b) )
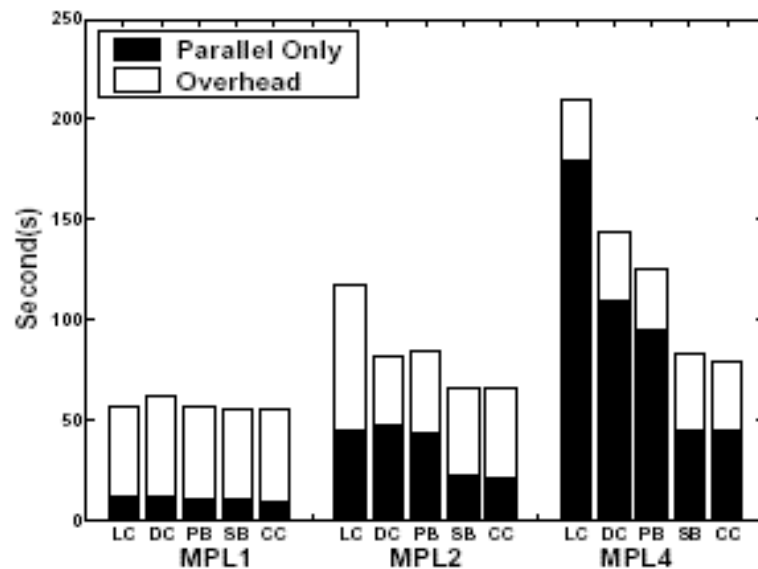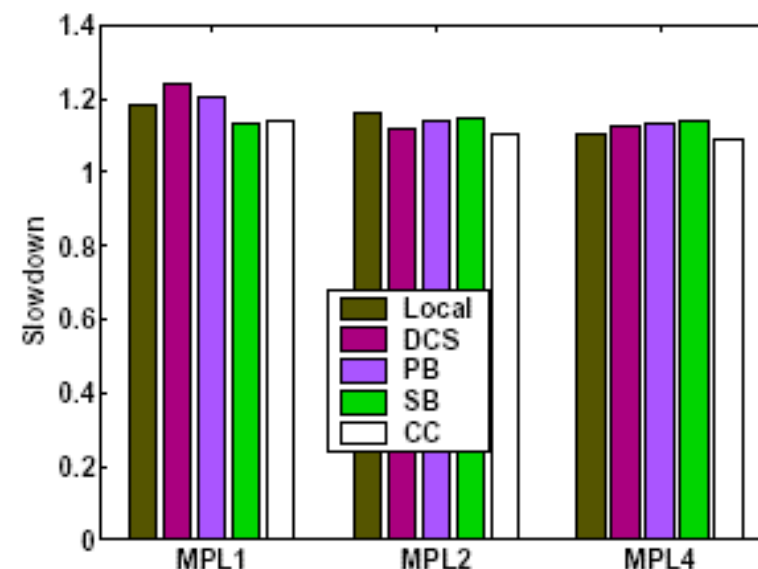


(a) Effect on parallel job

(b) Effect on CPU job

PENNSTATE

# Mixing I/O intensive jobs (W*I3*)

|| jobs in all schemes get equally affected (similar *overheads*).
Very Insignificant *slowdown* in I/O across all schemes.



(a) Effect on parallel job.          (b) Effect on I/O job.

# Conclusions

**Primary Contributions :-**

- Modular, flexible framework for deploying coscheduling.
- Fair, performing, new CC scheme.

**Significant findings :-**

- Blocking based schemes better than spinning (Linux).
- CC and SB scale well (Vertically) at high MPL of 6.
- CC and SB get equal or better than Batch scheduling.
- CC marginally better than SB, added QoS Potential.

**Other Advantages :-**

- Can implement all policies with CC approach.
- Can use framework with all ULN libraries.

PENNSTATE

# Future Work

- Horizontal scalability ( > nodes, GM)
- Remove assumptions :
  - No paging (All apps fit well in memory).
  - Homogeneous nodes availability.
  - All apps of same priority.
- Optimizations in CC mechanism.
- Allocation problems in coscheduling.
- Dynamic communication pattern identification.
- Integrated coscheduling as a feature in OS.
- True end-to-end QoS with support from scheduler.

# Questions ?

# Thanks for your Time !!

PENNSTATE

# fference with ICS ?
## (Optional slide)

- SR register for ALL incoming messages.
- ICS registers for expected messages for which a send has been done earlier.
- ICS is more tightly coupled (works on a send-recv pair).
- ICS is not too suitable in MPI environment.