

Job Communication Characterization and its Impact on Meta-scheduling Co-allocated Jobs in a Mini-grid

William M. Jones

Louis W. Pang

Dan Stanzione

Walter B. Ligon III

wjones@parl.clemson.edu

Parallel Architecture Research Laboratory

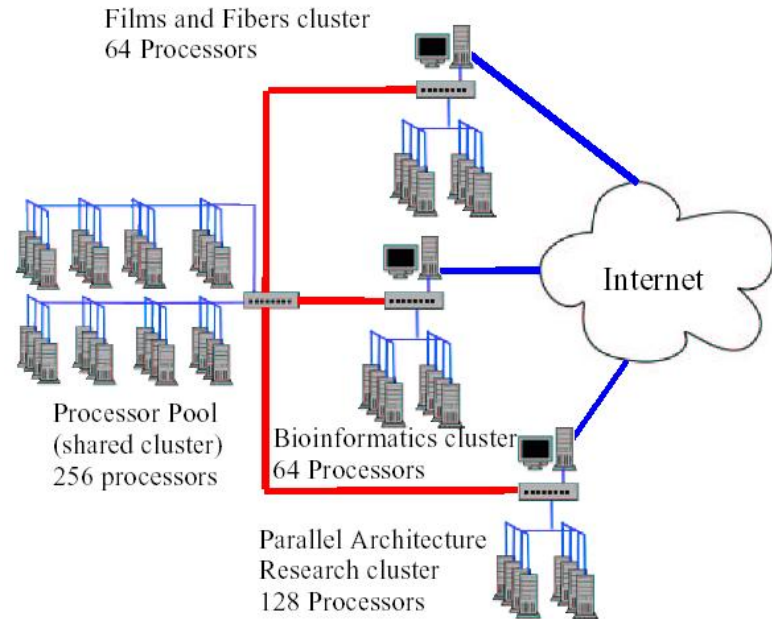
Clemson University

Clemson, South Carolina

<http://www.parl.clemson.edu>

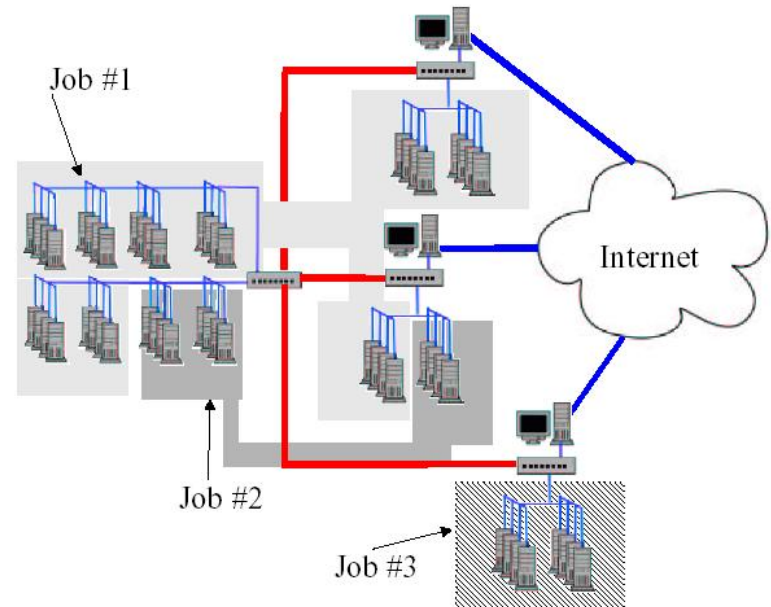
What is a Mini-grid ?

- Geographic co-location
- Dedicated network
- Fine-grain control
- Share resources



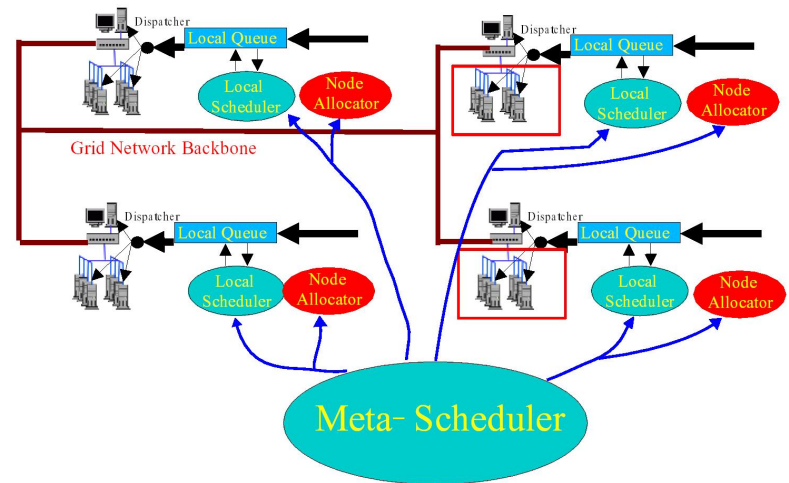
Job Scheduling

- Local execution
- Job migration
- Job co-allocation
 - ❑ Map across boundaries
 - ❑ Sharing resources
 - ❑ Network BW contention



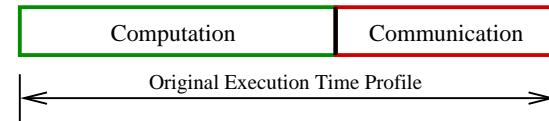
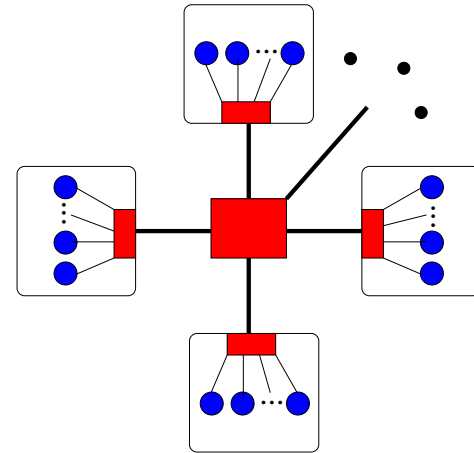
Initial Goals

- Meta-scheduler
- Co-allocation
- Create model
- Capture interaction
- Shared network
- Compare and contrast
- Improved scheduling



Models

- Mini-grid
 - ❑ Homogeneous nodes
 - ❑ Interconnection
- Parallel job
 - ❑ Comp. + Comm
 - ❑ Network congestion
- Communication
 - ❑ All-to-all
 - ❑ Per-processor BW

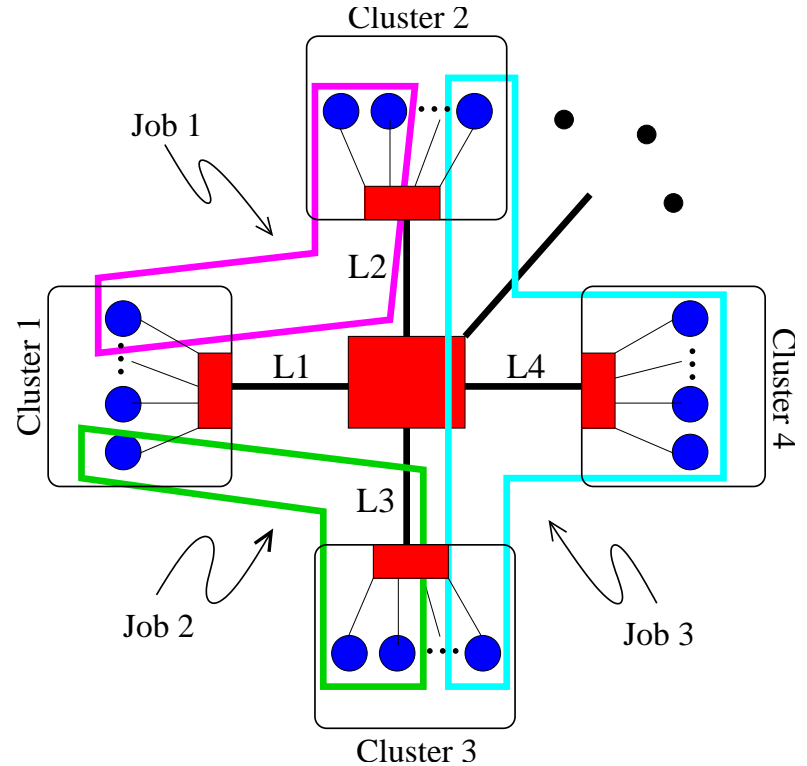
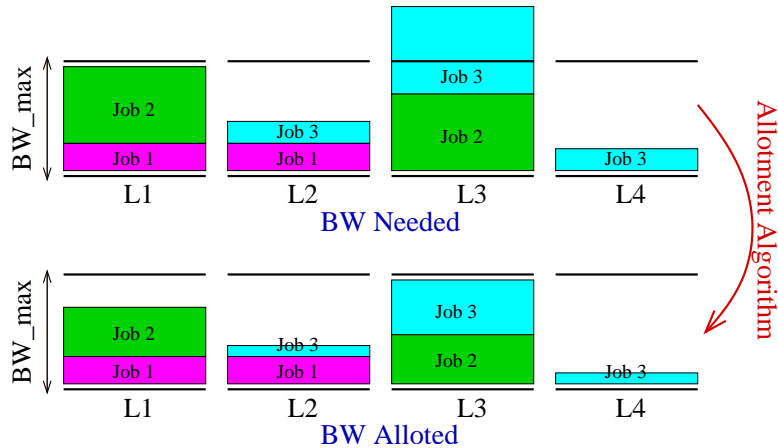


$$BW_i^j = \left(\frac{n_T^j - n_i^j}{n_T^j - 1} \right) (PPBW_j * n_i^j)$$

Job Co-allocation

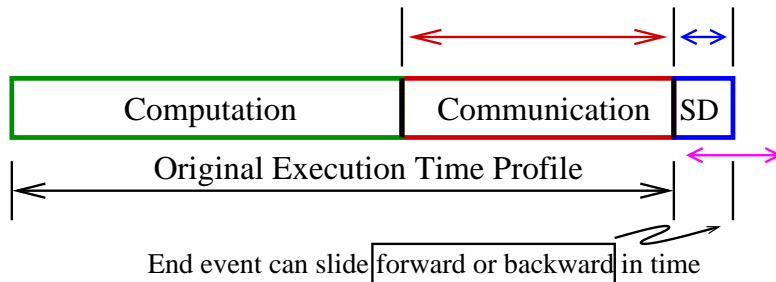
➤ Bandwidth saturation
 $BW_i^{sat} = BW_i^{max} / \sum_{\forall j \in J_i} BW_i^j$

➤ Communication slowdown
 $BW_{(k,j)}^{sd} = BW_{(k,j)}^{alloted} / BW_k^j$



Job Co-allocation Continued

- ⇒ State changing events
 - New co-allc'ed job
 - Co-allc'ed job finishes
- ⇒ Residual execution times



Residual Execution Time

$$T_E^R = \overbrace{T_C^R + T_P^R}$$

Residual Communication Time

$$T_C^R = \overbrace{(T_C^{R'} - T_C^\Delta)(BW_{(k,j)}^{sd'})} (BW_{(k,j)}^{sd})^{-1}$$

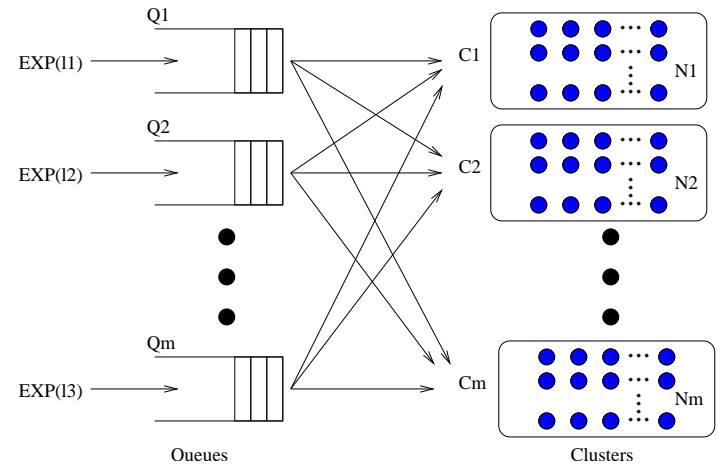
Residual Computation Time

$$T_P^R = \overbrace{(T_P^{R'} - T_P^\Delta)}$$

$$T_P^\Delta = \frac{\Delta T}{T_E^{R'}} T_P^{R'} \quad \text{and} \quad T_C^\Delta = \frac{\Delta T}{T_E^{R'}} T_C^{R'}$$

Experimental Setup

- Synthetic workload
- Poisson AP
- Exp. ST (*initially*)
- 70% comp. / 30% comm.
- 1 Gbps IC network links
- 4,000,000 jobs, $UNIF[10, 90]$
- 100 nodes per cluster



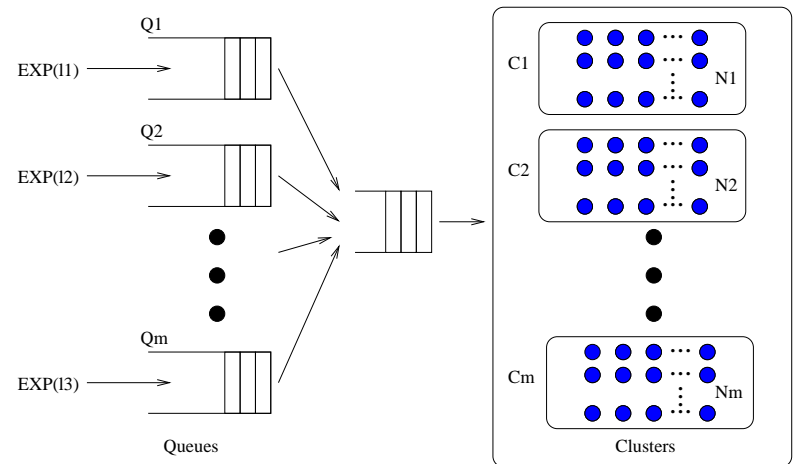
Simulation

➤ Algorithms

- No sharing
- No comm. penalty (Ideal)
- Migration only
- Fixed comm. model
- Dynamic comm. model

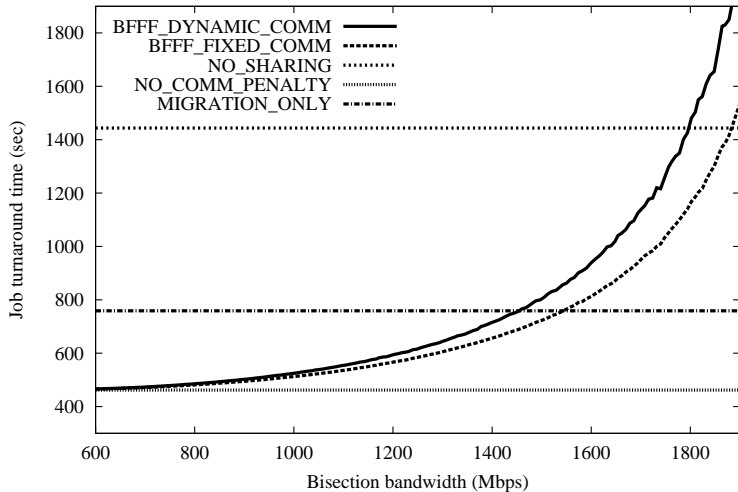
➤ Policies

- FCFS Scan
- Best Fit

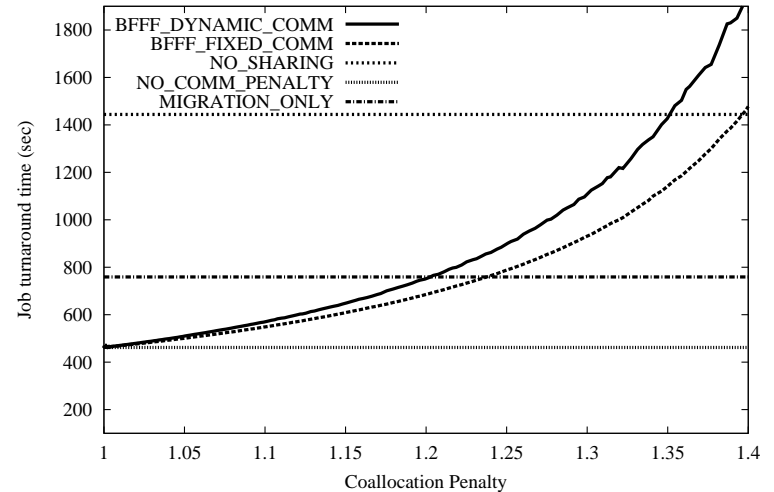


Results - 2 Cluster Case

Turnaround Time vs Bisection Bandwidth 4E6 Jobs / cluster, 2 clust. (100 nodes each)



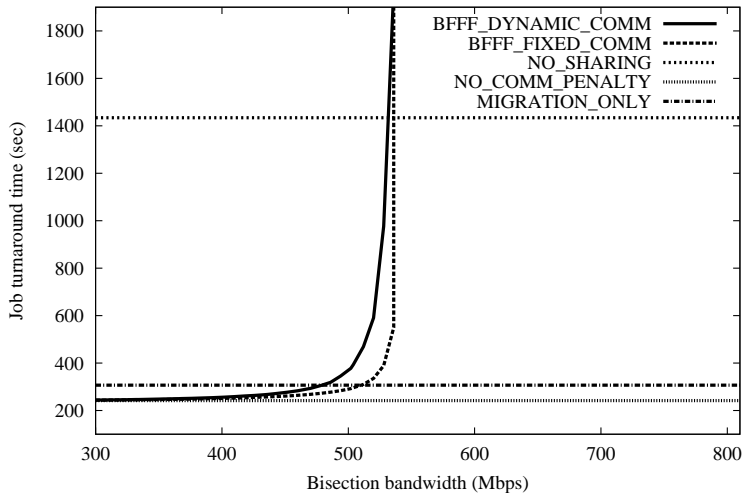
Turnaround Time vs Coalloc. Penalty 4E6 Jobs / cluster, 2 clust. (100 nodes each)



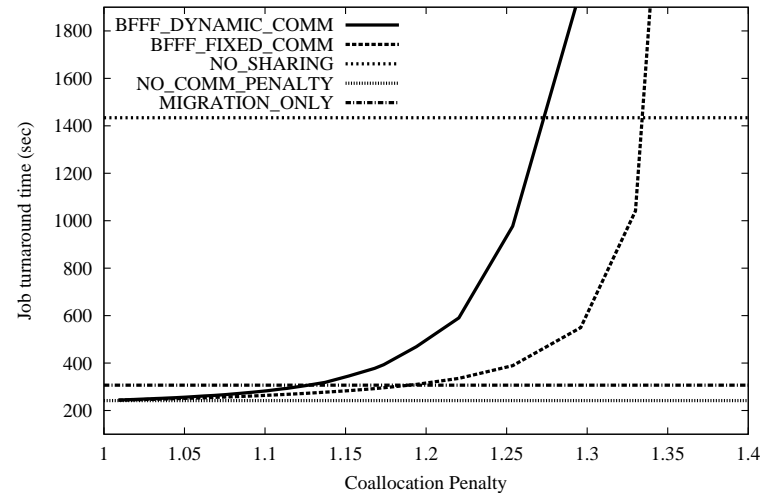
- Dynamic communication model is not as generous as fixed penalty model
- Effects become exacerbated as the number clusters is increased

Results - 8 Cluster Case

Turnaround Time vs Bisection Bandwidth 4E6 Jobs / cluster, 8 clust. (100 nodes each)



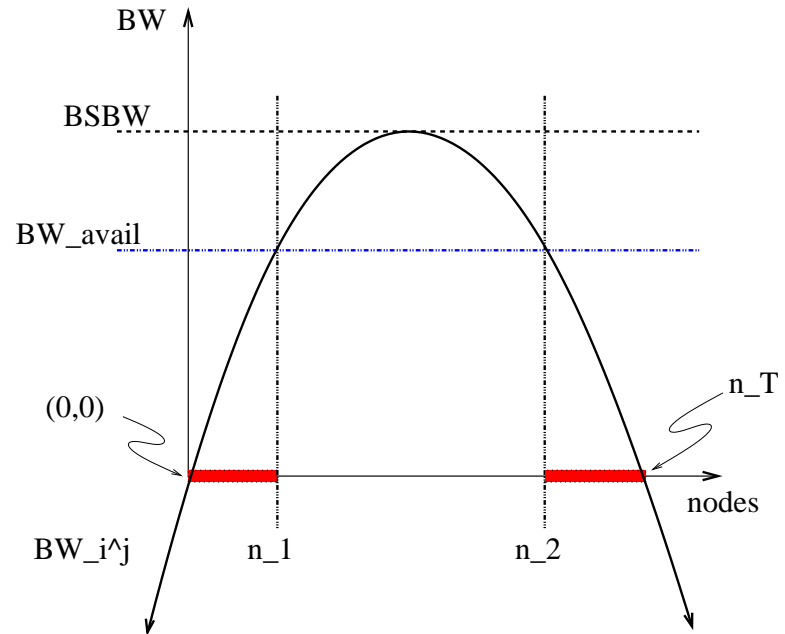
Turnaround Time vs Coalloc. Penalty 4E6 Jobs / cluster, 8 clust. (100 nodes each)

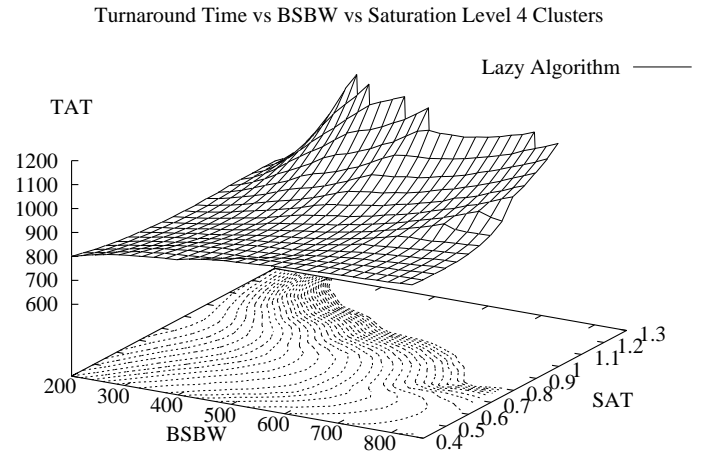
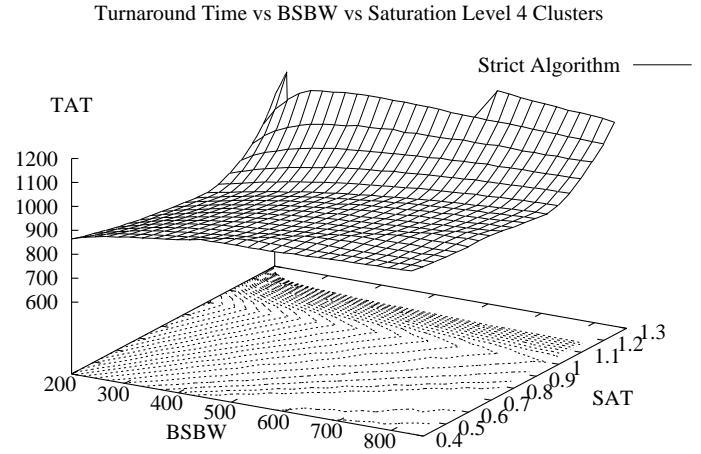
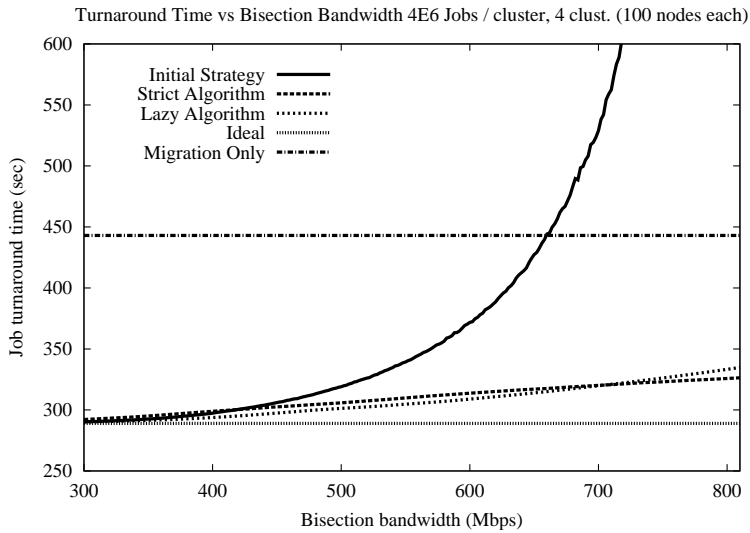


- Dynamic model accounts for the time varying utilization of intercluster NW links
- Captures the interaction of simultaneously co-allocated parallel jobs

Current Work

- ⇒ Constraint satisfaction
- ⇒ Optimization
- ⇒ New algorithms
- ⇒ Foreknowledge



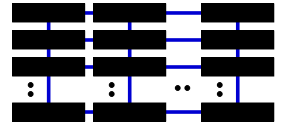


Future Work

- ⇒ Improve communication model
- ⇒ Verify saturation model with empirical data
- ⇒ Extend network to include additional topologies
- ⇒ Design more *intelligent* scheduling algorithms
- ⇒ Create more extensive priority mechanisms
- ⇒ Implement actual scheduler for our mini-grid

END

PARL
*Parallel Architecture
Research Laboratory*



Questions ?

Bandwidth Allotment Algorithm

$$BW_i^{uc_sat} = \frac{BW_i^{avail}}{\sum_{\forall j \in J_{uc}} BW_{(i,j)}^{alloted}}$$

- Step 1: Initialization** - For every job j , let $BW_{(i,j)}^{alloted} = BW_i^j$. For every link i , let $BW_i^{avail} = BW_i^{max}$. Let the unconstrained set of nodes, $J_{uc} = J$ (set of all jobs). Let the set, J_i be the set of all jobs that span link i .
- Step 2: Saturation detection** - For every link, calculate $BW_i^{uc_sat}$. While there exists at least one $BW_i^{uc_sat} < 1.0$, continue, else goto **Step 5**.
- Step 3: Saturation correction** - Identify the link with the smallest $BW_i^{uc_sat}$ (most saturated link) from **Step 2**, and globally reduce the allotted bandwidth of every job in $J_i \cap J_{uc}$ by a factor of BW_i^{sat} .
- Step 4: Update state** - Remove each of the modified jobs from the set J_{uc} . For each of the modified jobs, remove their allotted bandwidth from the available bandwidth, BW_i^{avail} on each link over which they span. Goto **Step 2**.
- Step 5: Termination** - DONE.