# Scheduling Many-Task Applications on Multi-Clouds and Hybrid Clouds

Shifat P. Mithila, Peter Franz, Gerald Baumgartner

# Who am I?

- Undergraduate research assistant
- Upgraded the platform to multiple cloud sites
- Learned algorithms through working with the code and explanations

# Status Quo Multi-Computer Approaches

| Approach | Performance | Networking | Resources | Parallelism |
|---|---|---|---|---|
| Supercomputers | Predictable | Low Latency + Robust | Homogeneous | Fine Grained Parallelism |
| Desktop Grids | Unpredictable | High Latency + Unreliable | Heterogeneous | No parallelism; Independent Tasks |
| Cloud Virtual Resources | Unpredictable | Moderate Latency + Relialiable | Heterogeneous | Coarse Parallelism; Map Reduce |

# Virtual Machine Challenges

- Performance information is unpredictable before resources are provisioned

- Latency to neighbors varies as host pcs change

- Heterogeneous resources

- Performance variability - competition for host os resources

# Supercomputer

- Forehand performance knowledge
- Robust and low latency networking
- Expensive
- Fine grained parallelism

# Large-scale desktop-based master-worker grids

- Heterogenous and variable performance of nodes
- Unreliable and high latency networking
- Donated CPU cycles
- Little to no parallelism

# Cloud Computing

- Virtual Machines
- Heterogenous and variable performance of nodes
- Robust but moderate latency networking
- Pay per use
- Fine grained parallelism (?)



There is no cloud
it's just someone else's computer

# Why you should be interested?

- Cheaper

- Higher availability

- Provider agnostic; allows price shopping

- Utilize your existing hardware

- Scale resources depending on task

A task based decentralized-vector-scheduling platform, minimizing propagation delay with latency based clustering.

# Task Based

- Tasks break large problems into smaller steps

- Some tasks require other tasks to be completed to start

- Tasks are python script to run with passed parameters



Spin-orbital coupled-cluster singles-and-double doubles equation
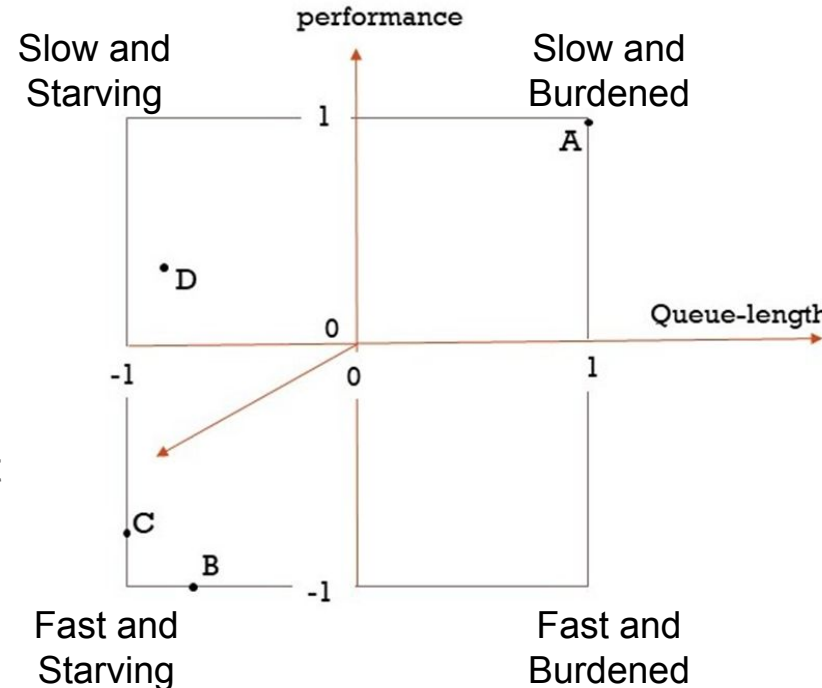
# Task Passing

- All nodes must be utilized

- Nodes must share

- Must be a fast decision

- Balancing act between multiple considerations
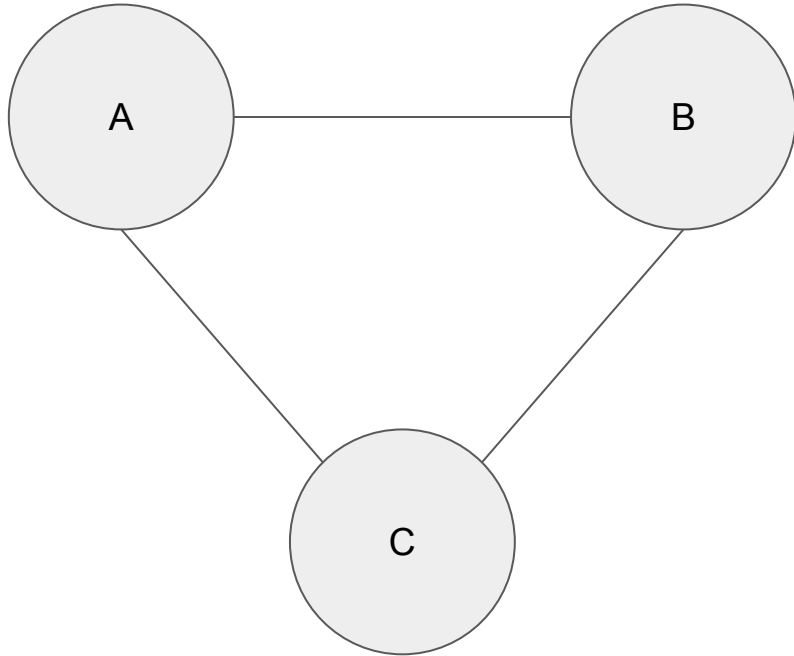
# Decentralized

- Avoiding a centralized bottleneck

- Short-running parallel-tasks on network connected heterogeneous resources

    - Controller is remote from nodes in cloud environment

    - Short running means less time for decisions to come from the controller

    - Controller needs to manage large amounts performance information

- Need task passing decision to be made at node level
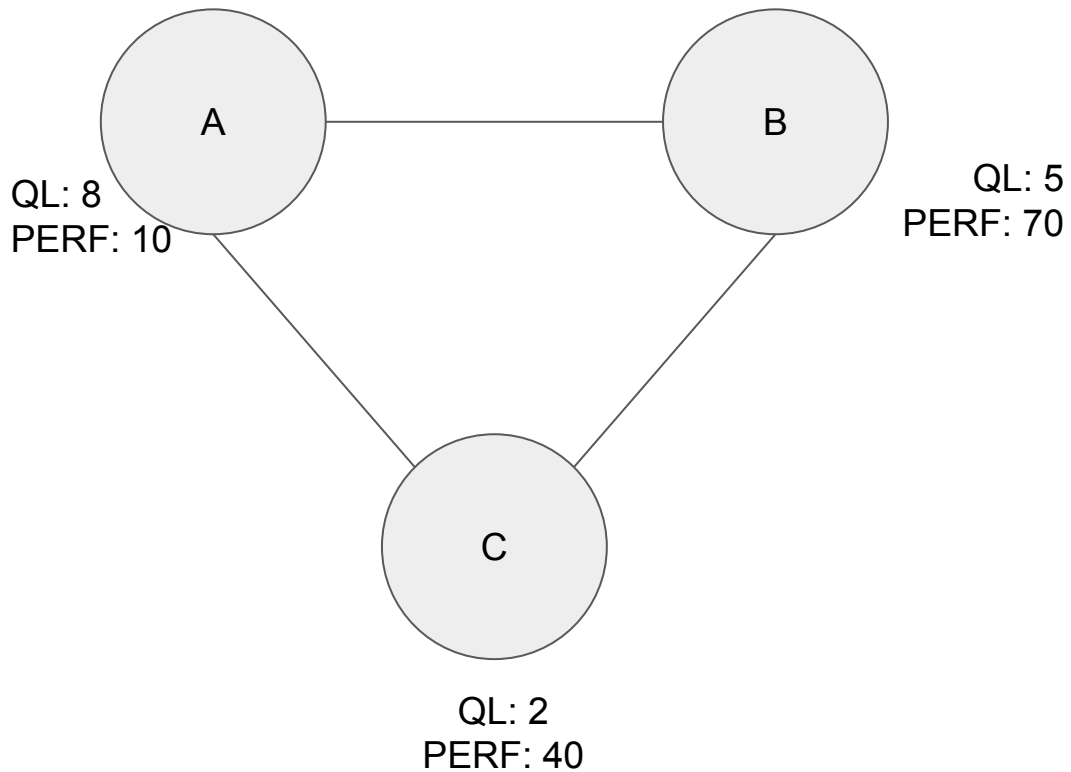
# Vector-Scheduling

- Lightweight decentralized algorithm

- Nodes connected in overlay graph

- Nodes send performance information to neighbors

- Measurements are normalized

- Using an experiment defined flow vector tasks sent in desired direction to neighbors



Slow and Starving

Slow and Burdened

Fast and Starving

Fast and Burdened

# Vector Scheduling worked example

A

QL: 8
Perf: 10

QL: 8
Perf: 10

B

QL: 8
Perf: 10

QL: 8
Perf: 10

C

Nodes provide information to neighbors

For each neighbor:
        shipQLToNeighbor(neighbor)
        shipPerfToNeighbor(neighbor)

A: [8,10]
B: [5,70]
C: [2,40]

A: [8,10]
B: [5,70]
C: [2,40]

A

B

C

A: [8,10]
B: [5,70]
C: [2,40]

Nodes keep track of the components in
a vector for each neighbor
[QueueLength, Performance]

```
For each neighbor:
    neighborVector =
    [
      reciveQL(neighbor),
      receivePerf(neighbor),
    ]
```

A: [1,-1]
B: [0,1]
C: [-1,0]

A: [1,-1]
B: [0,1]
C: [-1,0]

A

B

C

A: [1,-1]
B: [0,1]
C: [-1,0]

Nodes normalize vectors

# Scoring with a flow-vector

A

Flow Vector:
[-1,-0.25]

| Node | Queue Length | Performance | Score |
|------|--------------|-------------|-------|
| A | 1 | -1 | -0.75 |
| B | 0 | 1 | -0.25 |
| C | -1 | 0 | 1 |

A: (1 * -1) + (-1 * -0.25) = -0.75

B: (0 * -1) + (1 * -0.25)  = -0.25

C: (-1 * -1) + (0 * -0.25) = 1

# Thrashing

- Tasks passed back and forth with no completion
- Want passing, but not for minor gains
- "Statis vector" to minimize thrashing

# Stasis-Vector



A

Flow Vector:
[-1,-0.25]

Stasis Vector:
[0,-0.5]

| Node | Queue Length | Performance | Score |
|------|--------------|-------------|-------|
| A | 1 | -1 | ~~-0.75~~ -0.625 |
| B | 0 | 1 | -0.25 |
| C | -1 | 0 | 1 |

A: ( [1+0] * -1) + ([-1 + -0.5] * -0.25) = -0.625

# Adding latency

- Vector scheduling is an extremely expandable
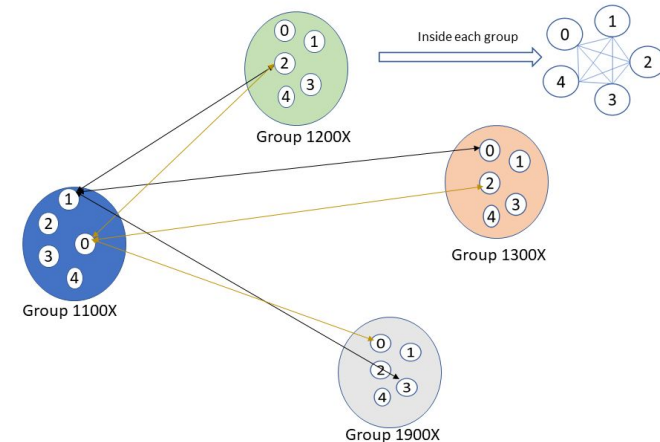
- Add Latency to scheduling

- Allows us to consider nodes relative locations

# Network Connectivity

- Full connectivity between all nodes has poor performance

- High latency connections should be avoided

- More neighbors increases overhead of vector scheduling

# Centralized Clustering Algorithm

- Measure latency to neighbors
- Nodes form a full graph with closest neighbors, Clique
- Keep some connections to distant groups


- Reduces number of neighbors to evaluate
- Removes slower connections
- Maintains connectivity

# Latency + Clustering increases performance

# Hybrid + Multi-Site Resource Proof of Concept
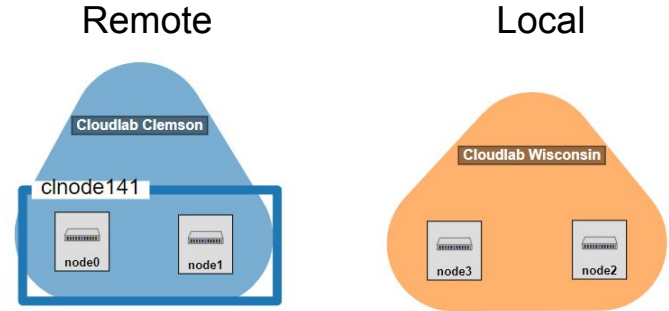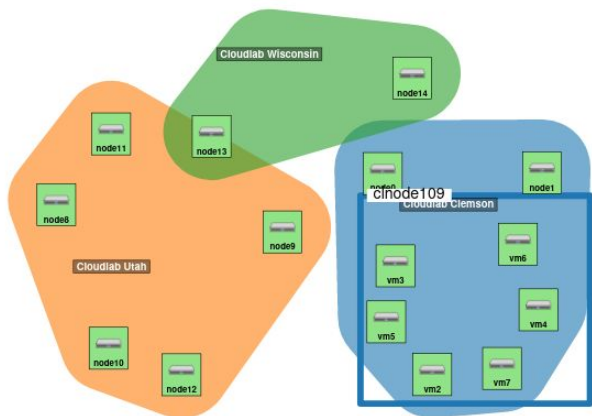
- 400 65 x 65 matrix multiplication tasks, with 50KB payload file included

- Demonstrates our capacity to stitch together multiple resource types and sites

- Limitation with handling LSU firewall, more development needed

| experiment | resources | initial vector | final vector | time |
|------------|-----------|----------------|--------------|------|
| Local | 2 PCs | [-1, -0.5, 0.5] | [-1, 0, -0.25] | 457.95 |
| Hybrid | 2 PCs, 2 VMS | [-1, -0.5, 0.5] | [-1, 0, -0.25] | 280.47 |



Remote

Cloudlab Clemson

clnode141

node0    node1

Local

Cloudlab Wisconsin

node3    node2

# Extending Research with Heterogeneity

- 400 65 x 65 matrix multiplication tasks, with 50KB payload file
- Heterogeneous performance
- Random graph construction, 20%



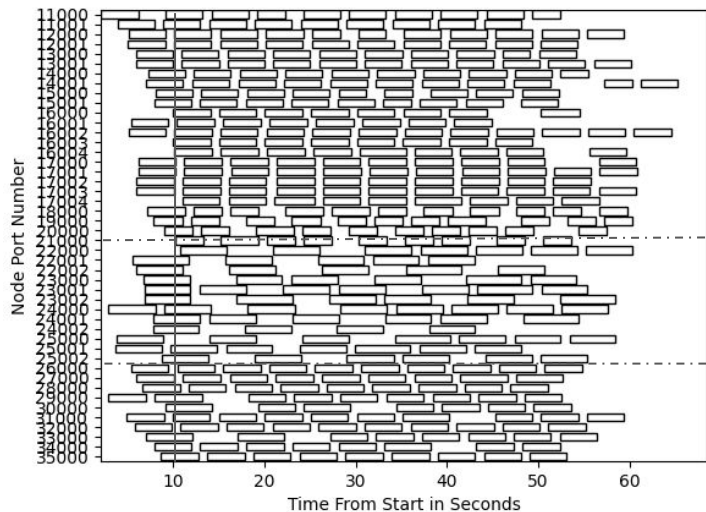| Site | Hardware Type | Workers | CPU | CPU Speed MHZ |
|------|---------------|---------|-----|---------------|
| Clemson | VM | 4 | 2 cores per | - |
| Clemson | c4130 | 10 | Intel E5-2680v3 | 2500 |
| Utah | d6515 | 10 | AMD EPYC 7452 | 2400 |
| Utah | m510 | 12 | Intel Xeon D-1548 | 2000 |
| Wisconsin | c220g2 | 10 | Haswell | 2600 |

# 3D Vector Scheduling in Multi-Site Cloud

- Varying vector configurations
- Ignoring latency had worst performance
- Avoiding latency is good for the middle and end
- Intentionally pushing tasks at the start, leads to good task saturation

Table 1: Performance comparison for different [queue, performance, latency] scheduling vectors with vector swap after 30 seconds.

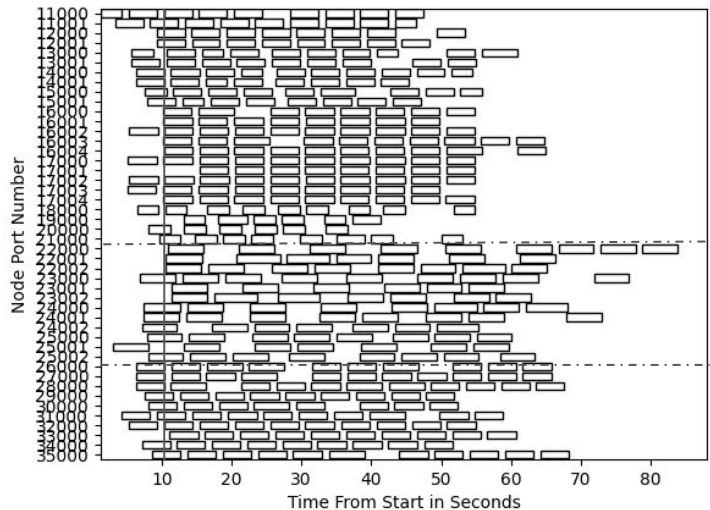| Experiment | initial vector | final vector | time |
|---|---|---|---|
| Ignore latency | [-1, -0.5, 0] | [-1, 0.5, 0] | 78.9 |
| Emphasize queue | [-1, 0, 0] | [-0.5, -0.5, -0.25] | 72.9 |
| Emphasize performance | [-1, -0.5, 0.5] | [-1, 0, -0.25] | 67.6 |

# Reading Tea Leaves

Use Latency + Push Away to Start



Ignore Latency



- Using latency to start has 93% Node utilization after 10 seconds
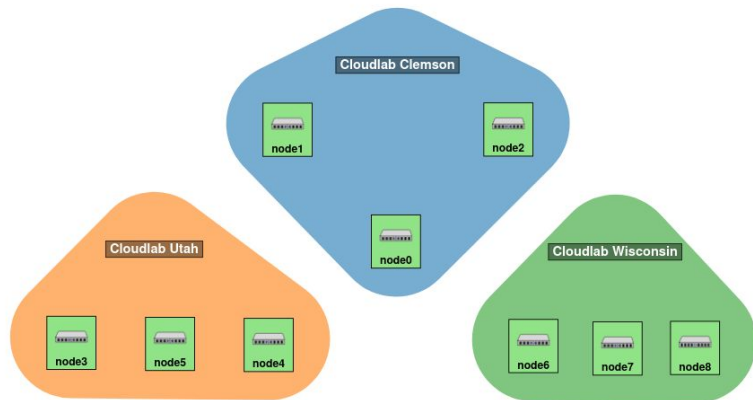- All sites have work until the end

- Ignoring latency has 73% Node utilization after 10 seconds
- Some sites starve at the end

# Extending Research with Heterogeneity

- 400 65 x 65 matrix multiplication tasks, with 50KB payload file

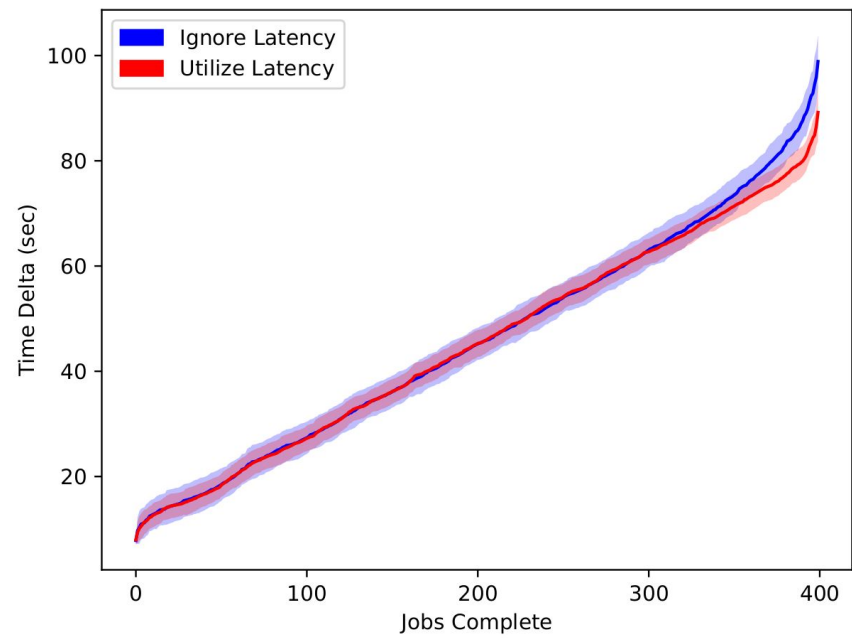- Heterogeneous performance

- Random graph construction, 20%



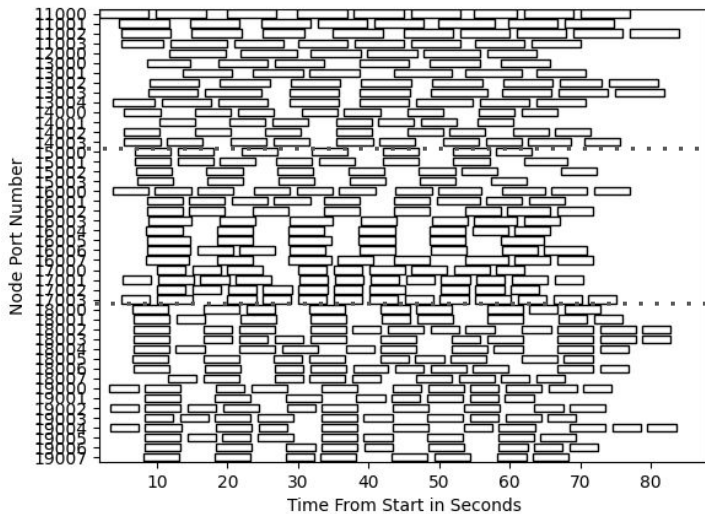| Site | Hardware Type | Workers | CPU | CPU Speed MHZ |
|---|---|---|---|---|
| Clemson | c8220 | 10 | Intel E5-2660v2 | 2200 |
| Utah | m510 | 16 | Intel Xeon D-1548 | 2000 |
| Wisconsin | c220g1 | 20 | Haswell | 2400 |

# 3D Vector Scheduling in Multi-Site Cloud

- Varying vector configurations

- Ignoring latency had worst performance

- Intentionally pushing tasks at the start, leads to good task saturation

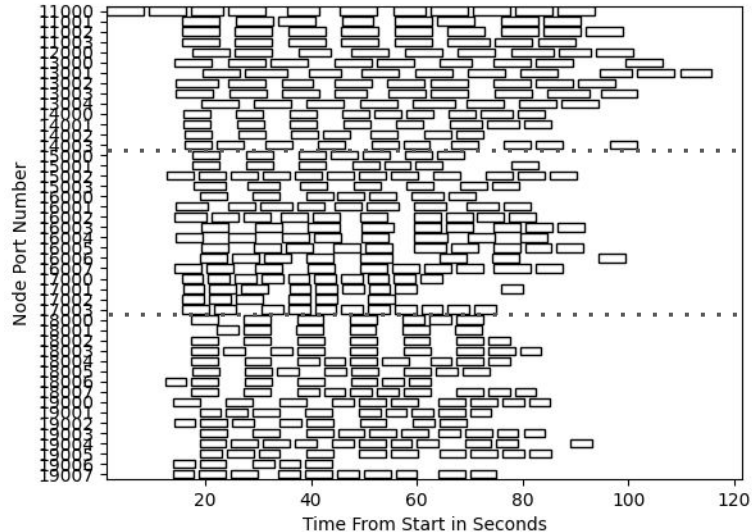| Name | Initial vector | Second Vector | Average Time |
|---|---|---|---|
| Ignore Latency | -1,-0.5,0 | -1,0.5,0 | 100.93 |
| Minimize Latency | -1,-0.25,-0.25 | -1,0,-0.25 | 91.55 |
| Push Tasks | -1,-0.5,0.5 | -1,0,-0.25 | 89.82 |

# Reading Tea Leaves



Use Latency + Push Away to Start

Ignore Latency
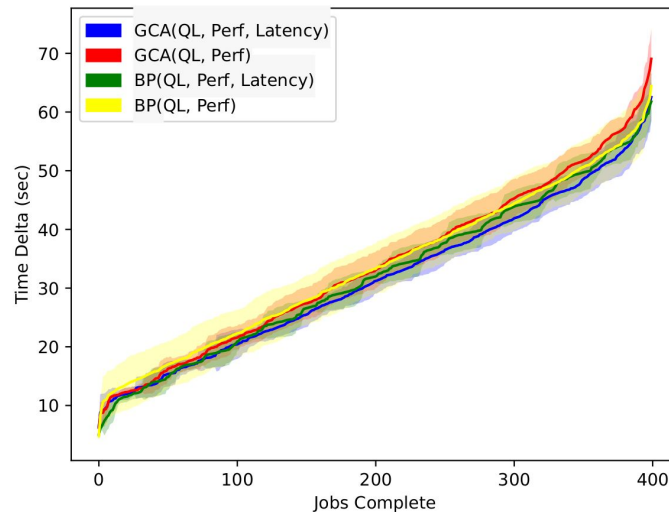
- Nodes quickly utilized
- All sites have work until the end

- Slow to utilize all nodes
- Some sites starve at the end

# Clustering Algorithm Limitations

- 46 worker nodes, 400 65x65 Matrix Mult tasks

  - GCA = Graph Construction Algorithm, Cluster

  - BP = Brain Peterson Algorithm, 20% Connectivity

- Clustering algorithm on real world clusters have oversize cliques

- Limited trials

# Further Work

- Improve adaptation of clustering algorithm parameters

- Trial automation


- General algorithm for changing flow vectors at runtime

- Probe hardware and system information to include as a vector component

- Task dependencies

- Tasks ran on collaboration of multiple nodes