



**Building the**

**International Data Placement Lab**

**Greg Thain**

**Center for High Throughput Computing**

# Overview

What is the IDPL?

How we built it

Examples of use

# Who is the IDPL?

Phil Papadapolus -- UCSD  
Miron Livny -- Wisconsin

Collaborators in  
China: Beihang // CNIC





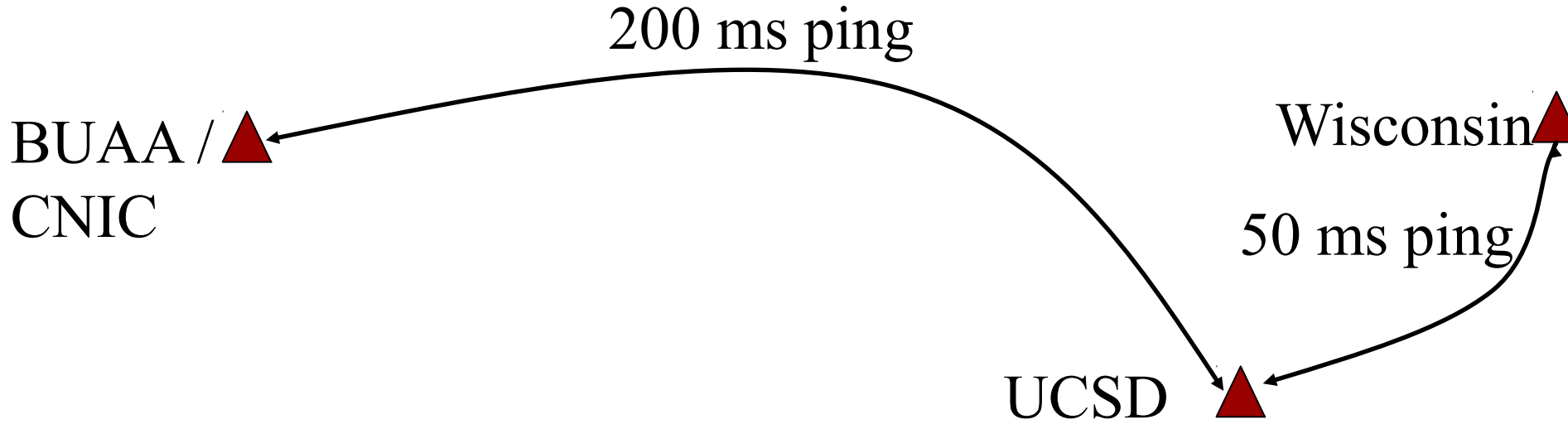
Beihang

CNIC

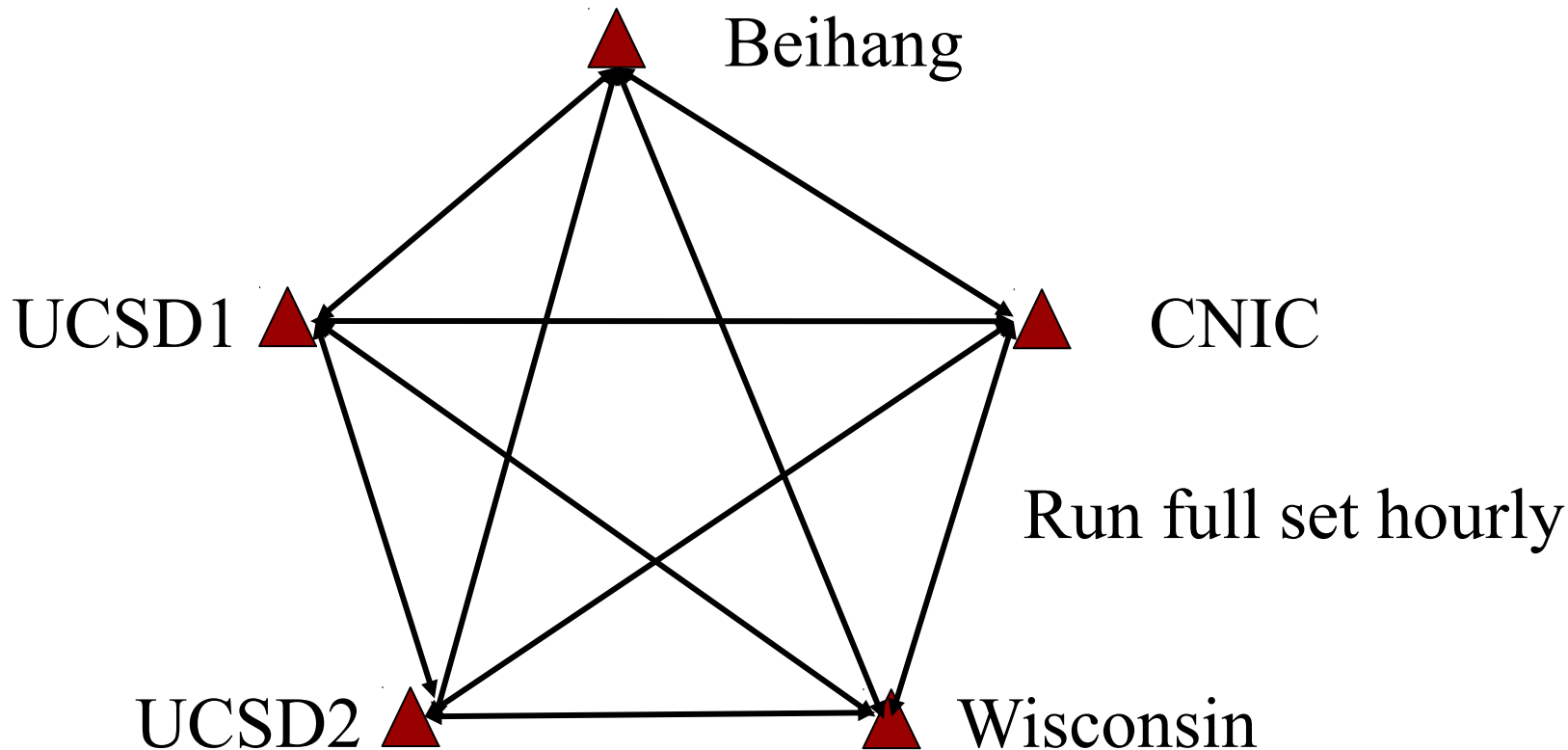
Wisc

UCSD (2x)

# A network engineer worldview:



# A Topologist's World View

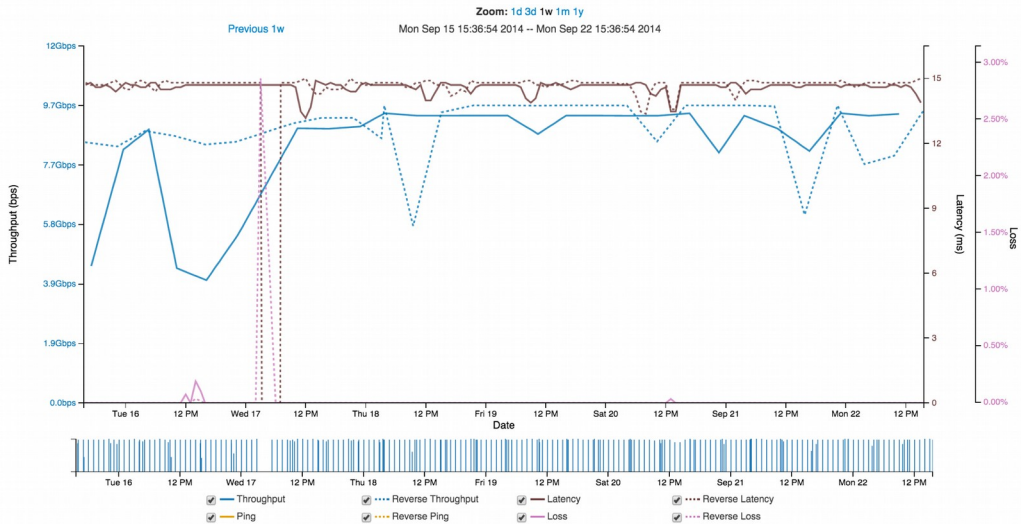


# Wait? What about PerfSonar?

# perfSONAR

[www.perfsonar.net](http://www.perfsonar.net)

Source	Destination	<a href="#">Link to this chart</a>
ani-pt1.es.net - 198.124.252.117 Capacity: 10G MTU: 9000	-- bnl-pt1.es.net - 198.124.238.38 [traceroute] Capacity: Unknown MTU: Unknown	
ani-pt1.es.net - 198.124.252.117 Capacity: 10G MTU: 9000	-- bnl-owamp.es.net - 198.124.238.49 Capacity: 10G MTU: 9000	
ani-owamp.es.net - 198.124.252.97 Capacity: 1.0G MTU: 1500	-- bnl-pt1.es.net - 198.124.238.38 Capacity: Unknown MTU: Unknown	
ani-owamp.es.net - 198.124.252.97 Capacity: 1.0G MTU: 1500	-- bnl-owamp.es.net - 198.124.238.49 [traceroute] Capacity: 10G MTU: 9000	



**NO PACKET LOSS AT EDGE  
ROUTER**

**THERE'S NO PROBLEM  
HERE**



# Wait? What about PerfSonar?

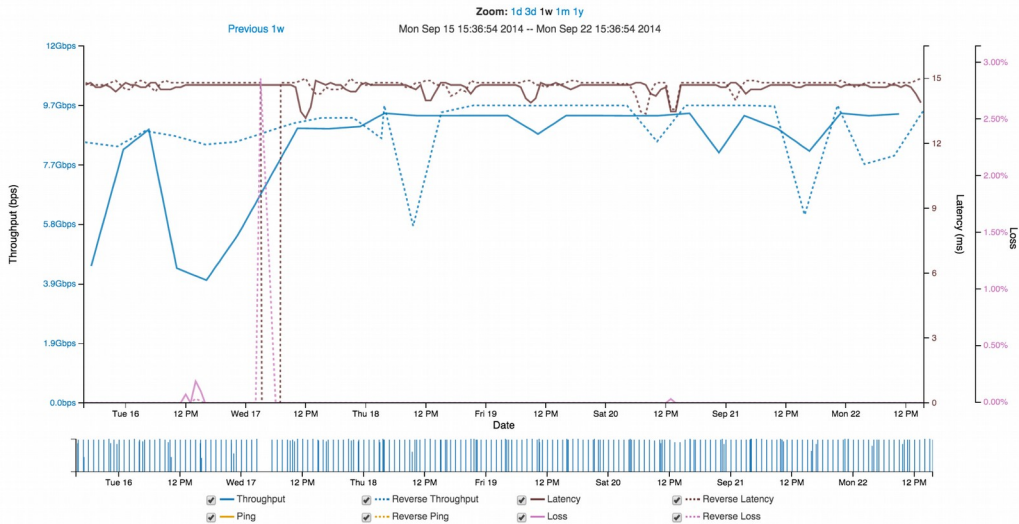
# perfSONAR

[www.perfsonar.net](http://www.perfsonar.net)

Necessary, not sufficient

Network only

Source	Destination
ani-pt1.es.net - 198.124.252.117 Capacity: 10G MTU: 9000	bni-pt1.es.net - 198.124.238.38 [traceroute] Capacity: Unknown MTU: Unknown
ani-pt1.es.net - 198.124.252.117 Capacity: 10G MTU: 9000	bni-owamp.es.net - 198.124.238.49 Capacity: 10G MTU: 9000
ani-owamp.es.net - 198.124.252.97 Capacity: 1.0G MTU: 1500	bni-pt1.es.net - 198.124.238.38 Capacity: Unknown MTU: Unknown
ani-owamp.es.net - 198.124.252.97 Capacity: 1.0G MTU: 1500	bni-owamp.es.net - 198.124.238.49 [traceroute] Capacity: 10G MTU: 9000



# Requirements for data placement

## 1. Co-scheduling of jobs

Start client and server simultaneously

## 2. Timed start of jobs

To prevent overlap and race conditions

## 3. Returning and managing of results

# Co-Scheduling with Parallel Universe

- Using static slots
- One slot per host
- Submit two-host job
  - Proc 0 is client
  - Proc 1 is server
    - (just a convention)

# Example Submit file

```
Universe = parallel
Executable = startup_script.sh
+ParallelShutdownPolicy = "WAIT_FOR_ALL"
<usual file transfer stuff>
Requirements = (machine == "client-machine-name")
Queue
Requirements = (machine == "server-machine-name")
Queue
```

# Startup Script (mark 1)

```
#!/bin/sh
if [ $_CONDOR_PROCNO = 0 ]
then
    do_client_stuff
else
    do_server_stuff
fi
```

# But this doesn't work...

- Synchronization problem
  - Not actually co-scheduled
  - Servers need to tell clients their port numbers

# Unix Process :: Condor Job

Unix Process	Condor Job
fork/exec	condor_submit
kill -9	condor_rm
Environment variables	Job Attributes
Standard error file	Condor job log

# condor\_chirp to the rescue!

- › condor\_chirp set\_job\_attr
- › Uses job ad as a ~~blackboard~~-whiteboard
- › Always read/writes to Proc 0  
(in parallel universe)



# Simplified Workflow

Server Side

Start Server

Listen on ephemeral port

```
condor_chirp set_job_attr Port #
```

Client Side

```
condor_chirp get_job_attr Port #
```

repeat as needed

```
run test
```

# That works – once

- › Need to do this periodically
- › Condor cron – two features

```
on_exit_remove = false
cron_minute = 23
cron_hour = 0-23/2
cron_month = *
cron_day_of_week = *
```

# Problems with condor cron

```
on_exit_remove = false
```

Means one job for all runs

stale set\_job\_attr's

held jobs a headache

Really want one condor job per run



Enter...

“CronMan”

# The Idea

- Have dagman itself be the restarter
  - Creates a new job every time
- With a one-node dag
- Whose one node has a delayed start time
  - Gives the placement job the job-nature

# The dag file

```
JOB A placement4-submit  
SCRIPT POST A /bin/false  
RETRY 1000000
```

# Modified Submit file

```
Universe = parallel
Executable = startup_script.sh
+ParallelShutdownPolicy = "WAIT_FOR_ALL"

cron_minute = 30

cron_window = 400

Requirements = (machine == "client-machine-name")
Queue
Requirements = (machine == "server-machine-name")
Queue
```

# condor\_q output

```
$ condor_q
```

```
ID          OWNER      SUBMITTED      RUN_TIME ST PRI  SIZE  CMD
26795.0     gthain     1/20 23:50     11+12:26:06 R  0    0.3   condor_dagman
27720.0     gthain     5/11 10:22       0+00:00:00 I  0    0.0   wrapper_script
27720.1     gthain     5/11 10:22       0+00:00:00 I  0    0.0   wrapper_script
```



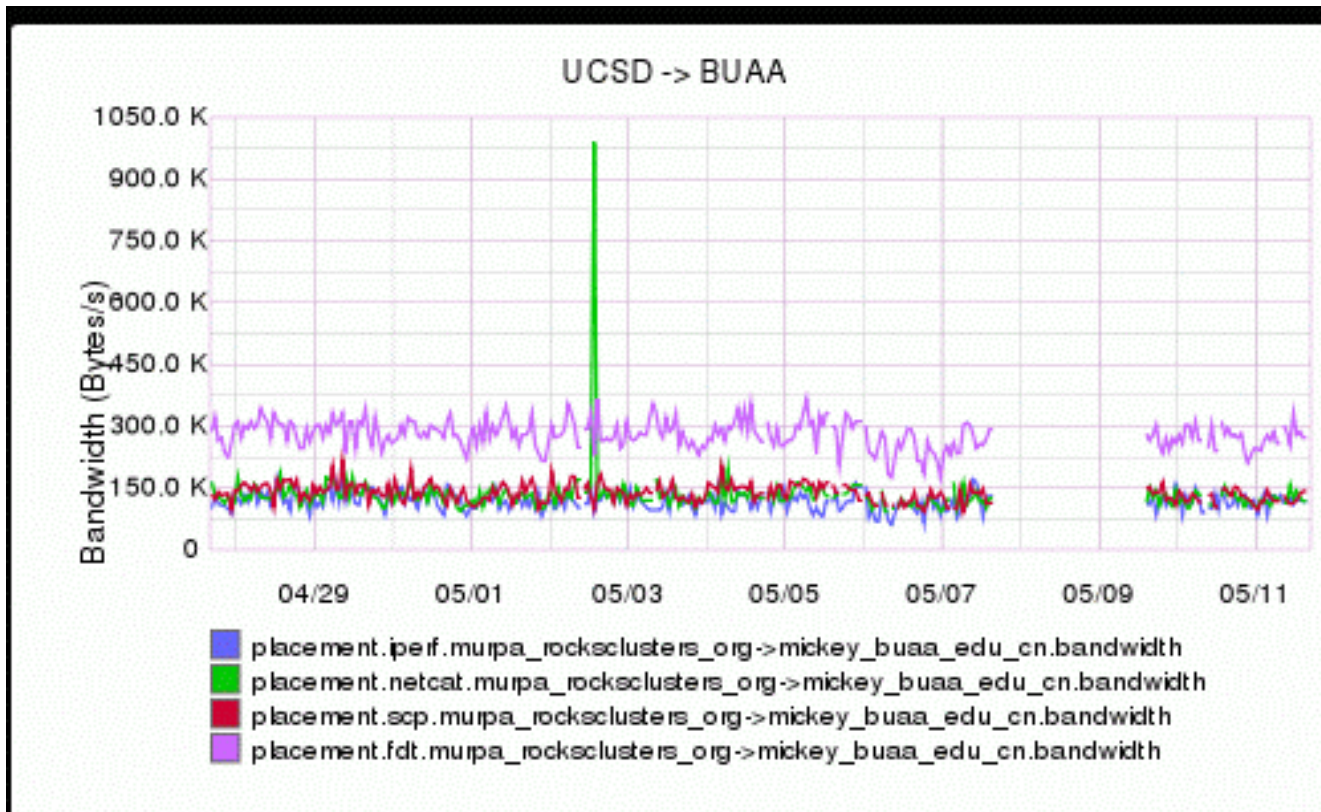
# Two of three problems solved

- What about reporting results?
- `condor_chirp ulog "server start"`
- `condor_chirp ulog "client results x y z"`

# Job log output

```
000 (27720.000.000) 05/11 10:22:48 Job submitted from host:      DAG Node: A
014 (27720.000.000) 05/11 11:17:00 Node 0 executing on host: client
008 (27720.000.000) 05/11 11:17:00 'client:start'
014 (27720.000.001) 05/11 11:17:04 Node 1 executing on host:
001 (27720.000.000) 05/11 11:17:04 Job executing on host: MPI_job
008 (27720.000.000) 05/11 11:17:06 'server:start'
033 (27720.000.000) 05/11 11:17:08 Setting job attribute iperfServer to '20650'
008 (27720.000.000) 05/11 11:17:31
'results,1462979830.739976,1462979851.114580,1,20.100000,5736960'
008 (27720.000.000) 05/11 11:17:31 'komatsu.chtc.wisc.edu(iperf) client:end'
008 (27720.000.000) 05/11 11:17:31 'komatsu.chtc.wisc.edu(irods) client:start'
008 (27720.000.000) 05/11 11:17:31 'davos.cyverse.org(iperf) server:end'
```

# Then forward to Graphite



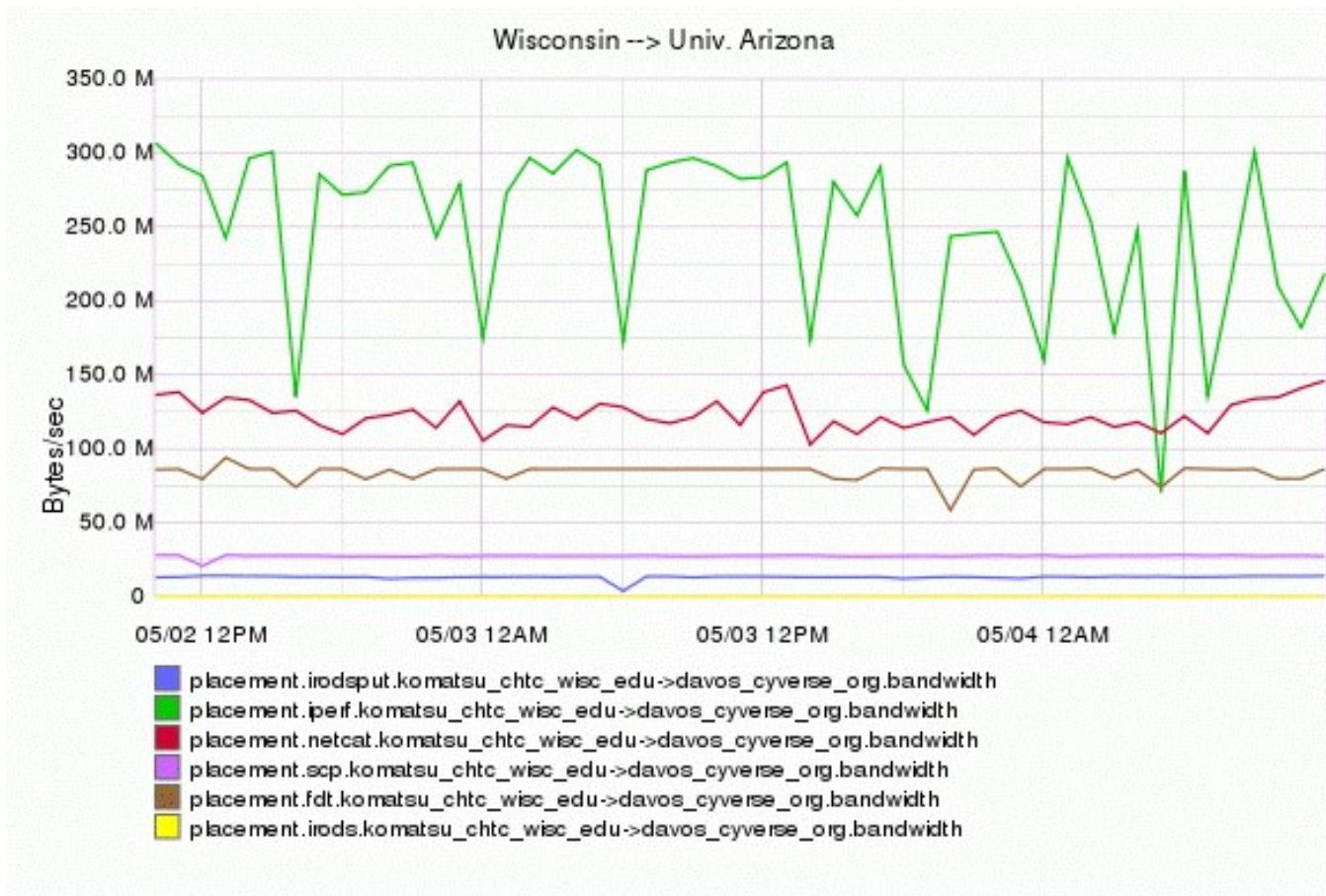
# Application of IDPL: The Phytomorph problem

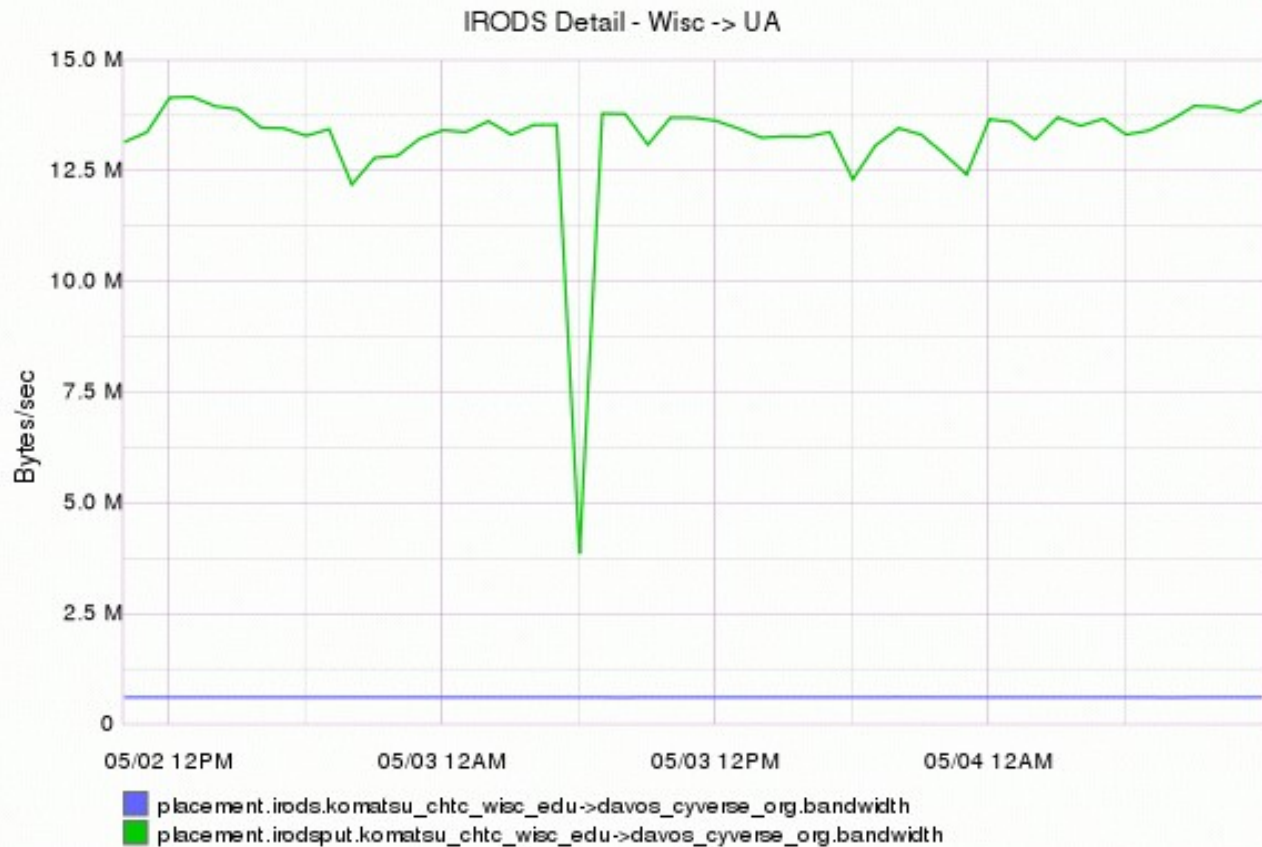
PI in the Spalding lab researching corn seedlings

Runs analysis jobs at Wisconsin HTCondor pool

Pulling data from Cyverse (nee iPlant) in Arizona

Claims data xfer rates slow



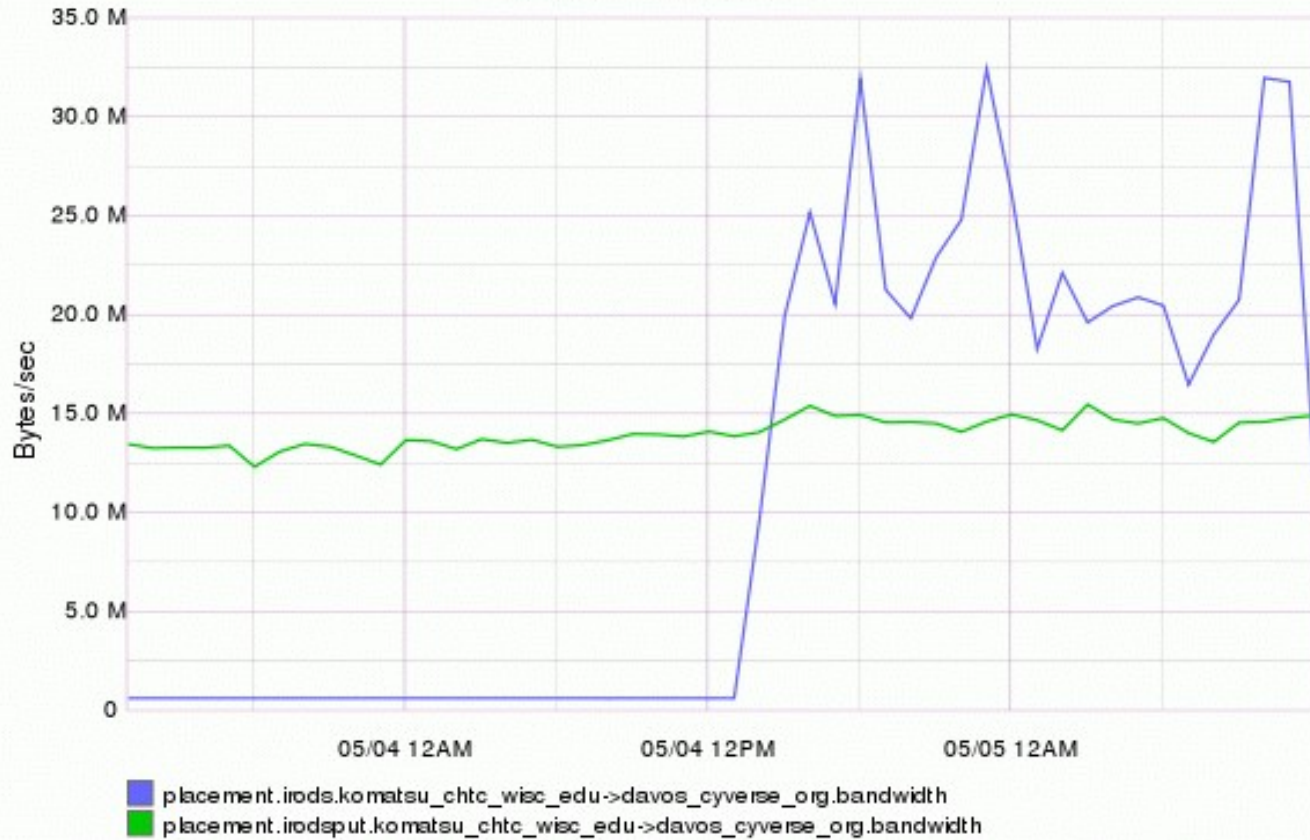


# After talking with iRODS devs

- All performance testing had been LAN
- Our Networking folks ID'd problems
- Made one big fix
  - Don't set TCP\_SNDBUF explicitly!
- Gave us new clients & servers



IRODS Detail - Wisc -> UA





# Future Work

Add more protocols

More sites

Work with clouds

# Thank you!

- Try CronMan pattern yourself!
  - Even with vanilla universe jobs
- Think about Unix process patterns
- Example on slides simplified!
- Real code on git hub at
  - <https://github.com/iDPL/placement>