

HTCondor workflows at Utility Supercomputing Scale: How?

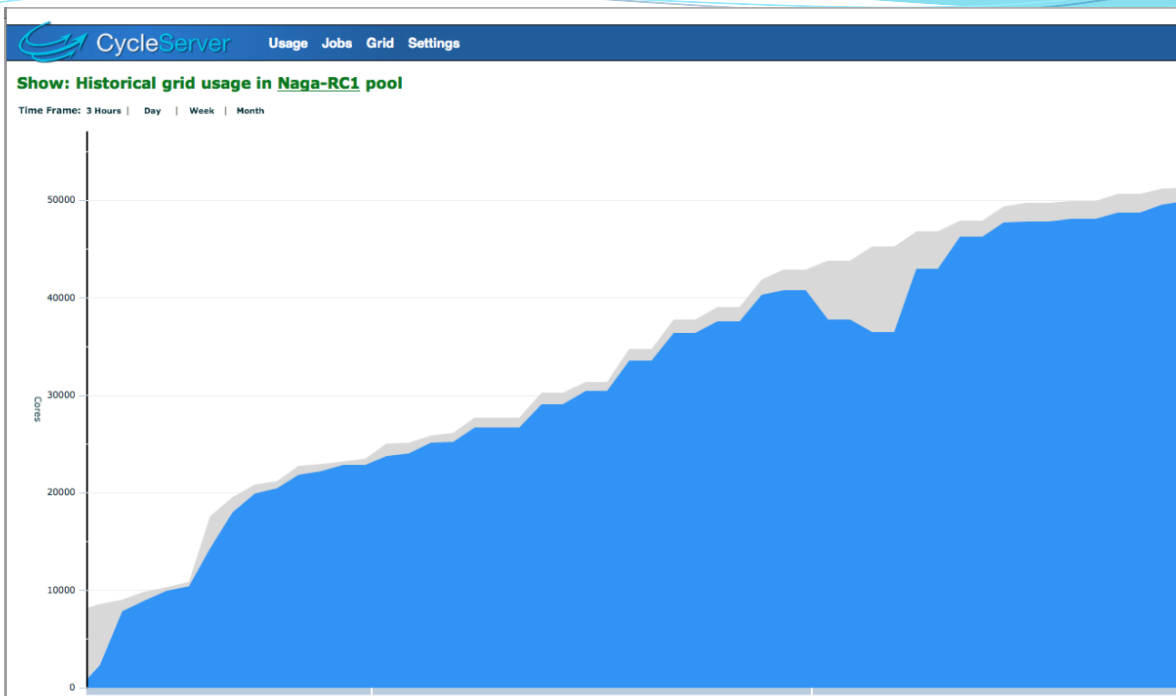
Ian D. Alderman
Cycle Computing

Thundering Herd Problem



Thundering Herd Problem

- Classical OS problem: multiple processes are waiting for the same event, but only one can respond at a time.
- In the cloud, what happens to the (underlying) infrastructure when you start 10k servers is someone else's problem.
- What happens at the platform and application level is your problem
- Experience is helpful.



Ramping up to
50,000 cores

**I DON'T OFTEN RUN 50,000
CORES**



memegenerator.net

while true bottleneck.next()

- From Miron:
 - **A bottleneck is a (system) property that once removed creates a new bottleneck.**
- Related to theory of constraints from industrial engineering.
- Corollary: Every component in a distributed system can be a bottleneck.

Bottlenecks we have seen

- Scheduler. Forking, transferring data, etc.
- Shared filesystem (NFS).
- Web server/backend/provisioning system – client.
- Provisioning system - server (AWS). Need delta mechanism for ec2-describe-instances.
- Configuration management system. Designed to handle updates in large systems, not provision large systems all at once.



Message in a bottleneck?

Find the right problem: Aim high.

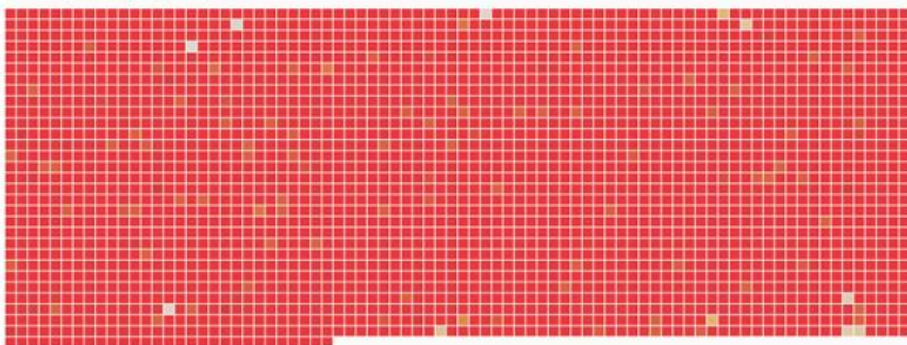
- Predict costs, runtime. Understand I/O and memory requirements. Users don't always know this.
- Zach says: **Understand your job**. Users don't often have the tools to do this.
- We were surprised to find out that Flexera license server can handle this scale given enough file handles.
- **The right bottleneck is CPU: that's what we're paying for**

Distributing jobs

- Distribute tasks among several scheds. (Manure spreaders)
- CycleServer manages tasks across several environments.
- Multi-region, heterogeneous clusters.
- Goals:
 - Keep queues filled (but not too full)
 - Keep queues balanced
 - Minimize complexity
 - Reduce server overhead costs

Cluster summary for all clusters

Show: CPU Usage by Host



View: Performance Download CSV Refresh Show Detail View Detail Search:

Host	Cluster	Clock Speed	Memory	CPU Usage	Mem Usage	Net In	Net Out	Last Reported
internal		2.00 GHz	15.00 GB	100%	5%	33.50 B/s	751.38 B/s	7:17 PM
internal		2.00 GHz	15.00 GB	100%	5%	38.53 B/s	763.60 B/s	7:11 PM
2.internal		2.00 GHz	15.00 GB	100%	5%	41.17 B/s	1.28 kB/s	7:17 PM
internal		2.00 GHz	15.00 GB	100%	5%	33.78 B/s	1.39 kB/s	7:18 PM
internal		2.27 GHz	7.50 GB	100%	6%	77.70 B/s	624.97 B/s	7:12 PM

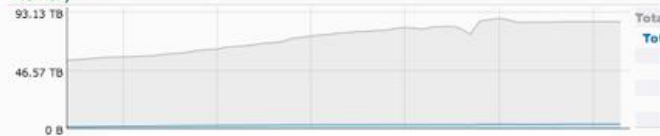
Cluster Performance Stats

Time Frame: Hour

CPU



Memory



Network



CycleCloud: Auto-start and auto-stop at the cluster level

- Automation is the goal: nodes start when jobs are present, nodes stop when jobs aren't there (5 minutes before the billing hour mark).
- Select instance types to start in rank order to maximize price-performance.
- Use pre-set spot prices to minimize costs.

Zero-impact job wrapper

Goal: Don't hit the file server, don't have HTCondor transfer anything.

- No file transfer
- No input
- No results
- No output, error or log
- So how does the job do anything?

Use S3 instead of file server

- B3: bottomless bit bucket.
- Eventual consistency is well suited for the type of access patterns we use:
 - Read (big) shared data
 - Read job-specific data
 - Write job-specific results
 - Jobs can be made to except (hold) when inputs aren't available (rare)
- Some systems do scale; this is one.

```

3
4 opts = GetoptLong.new( [ '--verbose', '-v', GetoptLong::NO_ARGUMENT ],
5                       [ '--overwrite', '-o', GetoptLong::NO_ARGUMENT ],
6                       [ '--stdout', '-s', GetoptLong::NO_ARGUMENT ],
7                       [ '--config', '-c', GetoptLong::REQUIRED_ARGUMENT] )
8
9 param = Hash.new
10 param[:verbose] = false
11 param[:stdout] = false
12 param[:overwrite] = false
13
14 opts.each do |opt, arg|
15   case opt
16   when '--verbose'
17     param[:verbose] = true
18   when "--stdout"
19     param[:stdout] = true
20   when "--overwrite"
21     param[:overwrite] = true
22   when "--config"
23     param[:config] = arg
24   end
25 end
26
27 class GlideJobWrapper
28   def initialize(a, p)
29     @verbose = p[:verbose]
30     @stdout = p[:stdout]
31     @overwrite = p[:overwrite]
32
33     @config = /

```

Don't overwrite results

```
def output_file_exists?  
  @log.debug "Checking for output file at #{output_file_url}."  
  `#{@s3cmd} ls #{output_file_url} | wc -l`.chomp.to_i > 0  
end
```

Actual check to see if
results are there already


```
def exp_backoff_retry_command(cmd, max=15, sleep=0.5, timeout=10*60)
  start_time = Time.now.to_i
  end_time= start_time+timeout
  count = 0
  @log.info "Attempting cmd: '#{cmd}'."
  while not system(cmd)
    exit_status = $?[0]
    if Time.now.to_i > end_time
      raise "Timeout #{timeout}s exceeded with cmd '#{cmd}'."
    end
    count = count + 1
    if count > max
      raise "Count max #{max} exceeded with cmd '#{cmd}'."
    end
    @log.error "Command failed, #{exit_status}, retrying..."
    sleep_time = sleep * 2**count
    sleep rand(sleep_time)
  end
  @log.info "Command succeeded."
end
```

Exponential back-off for data transfer

```
def glide_cli
  "#{@glide} -NOJOBID #{in_file} > #{in_file}.stdout 2> #{in_file}.stderr"
end

def my_file_p3(f)
```

Actual command
line captures
stdout and stderr

```
def mv_all_logs
  logs = %w(out stdout stderr)
  logs.each do |l|
    mv_file_s3 "#{in_file}.#{l}"
  end
end
```

If command succeeds,
save stdout and stderr

Actual submit file

universe = vanilla

Requirements = (Arch =?= "X86_64") && (OpSys =?= "LINUX")

executable = /ramdisk/glide_job_wrapper.rb

should_transfer_files = if_needed

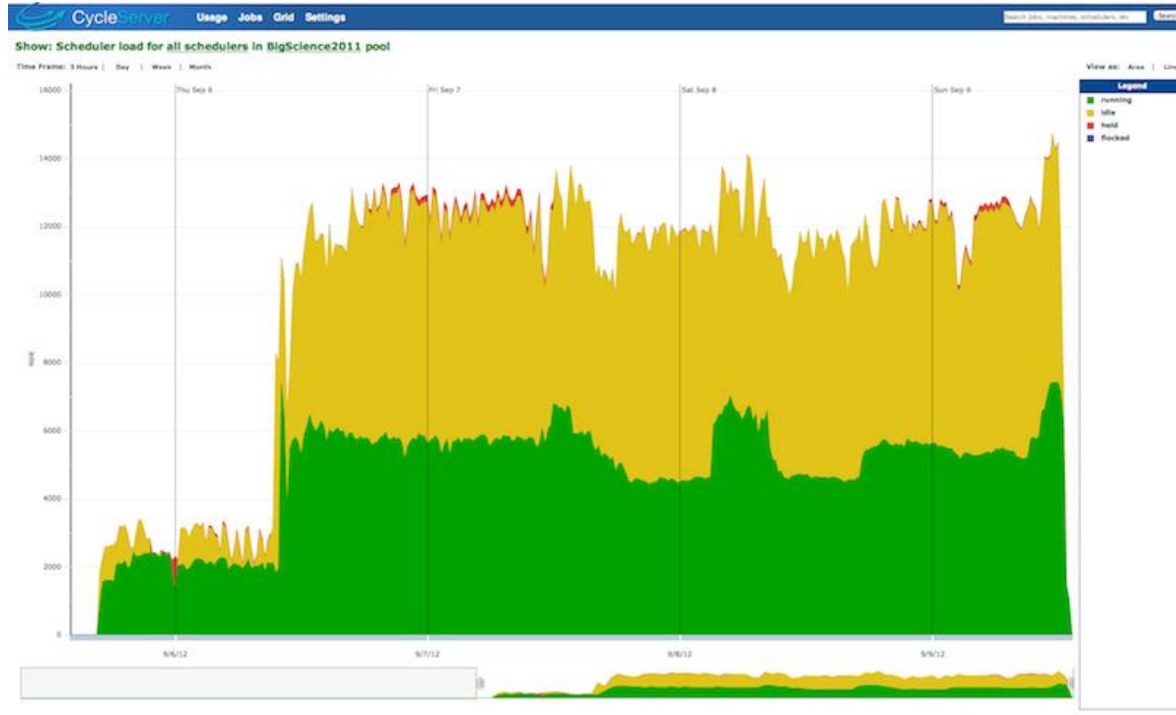
when_to_transfer_output = on_exit

environment = "..."

leave_in_queue = false

arguments = \$(process)

queue 325937



DAGMan is your friend

Configuration management system

- OpsCode Chef.
- Chef-solo.
- Chef Server 11 from OpsCode.
- Deploy changes to wrapper scripts, HTCondor configuration, etc during a run.
- Run OOB task on all hosts (knife ssh). Very cool but realistically can be a bottleneck.

Chef overview for chef-server-11.cycloid.com

Current host stats

Chef Servers: 1
 # Hosts: 10343
 # Converged Hosts: 10312
 # Unconverged Hosts: 31

Converge stats (last hour)

Total Converges: 3944
 Successful Converges: 3852
 Failed Converges: 92

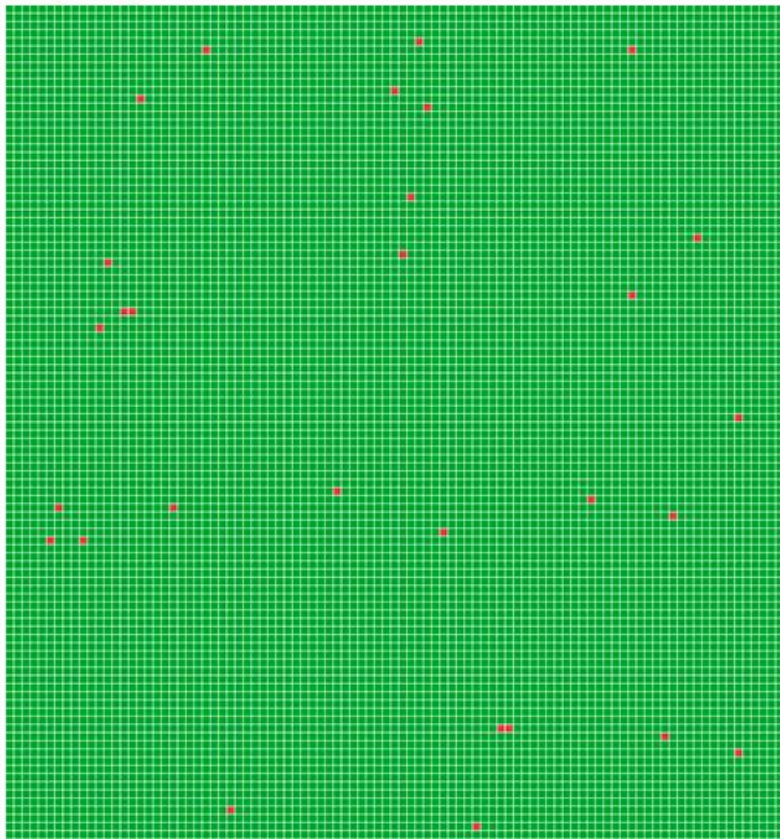
Alerts

Fri Feb 01 2013 18:21:53 GMT-0500 (EST): ec2-23-22-131-239.compute-1.amazonaws.com failed to converge

Recent converges

Time	Host Name	Status	Start Time	End Time	Duration
1:04 AM	ec2-54-234-122-169.compute-1.amazonaws.com	Success	1:04 AM	1:04 AM	12m 6s
1:04 AM	ec2-23-20-172-129.compute-1.amazonaws.com	Success	1:04 AM	1:04 AM	6m 19s
1:04 AM	ec2-54-242-78-85.compute-1.amazonaws.com	Success	1:04 AM	1:04 AM	6m 40s
1:04 AM	ec2-54-242-93-226.compute-1.amazonaws.com	Success	1:04 AM	1:04 AM	6m 22s
1:04 AM	ec2-184-73-138-253.compute-1.amazonaws.com	Success	1:04 AM	1:04 AM	6m 36s
1:04 AM	ec2-54-242-240-184.compute-1.amazonaws.com	Success	1:04 AM	1:04 AM	6m 17s
1:04 AM	ec2-107-20-144-205.compute-1.amazonaws.com	Success	1:04 AM	1:04 AM	6m 20s
1:04 AM	ec2-23-20-234-71.compute-1.amazonaws.com	Success	1:04 AM	1:04 AM	6m 9s
1:04 AM	ec2-54-242-80-75.compute-1.amazonaws.com	Success	1:04 AM	1:04 AM	6m 15
1:04 AM	ec2-50-16-169-53.compute-1.amazonaws.com	Success	1:04 AM	1:04 AM	6m 46
1:04 AM	ec2-54-242-44-228.compute-1.amazonaws.com	Success	1:04 AM	1:04 AM	6m 13s
1:04 AM	ec2-54-242-44-228.compute-1.amazonaws.com	Success	1:04 AM	1:04 AM	6m 20s
1:04 AM	ec2-54-242-59-145.compute-1.amazonaws.com	Success	1:04 AM	1:04 AM	6m 46
1:04 AM	ec2-54-242-202-232.compute-1.amazonaws.com	Success	1:04 AM	1:04 AM	6m 20s
1:05 AM	ec2-67-202-71-38.compute-1.amazonaws.com	Success	1:05 AM	1:05 AM	6m 46
1:05 AM	ec2-174-129-61-233.compute-1.amazonaws.com	Success	1:05 AM	1:05 AM	5m 55s
1:05 AM	ec2-174-129-137-142.compute-1.amazonaws.com	Success	1:05 AM	1:05 AM	5m 47s
1:05 AM	ec2-54-243-14-15.compute-1.amazonaws.com	Success	1:05 AM	1:05 AM	6m 10s
1:05 AM	ec2-23-22-137-89.compute-1.amazonaws.com	Success	1:05 AM	1:05 AM	6m 2s
1:05 AM	ec2-54-242-240-75.compute-1.amazonaws.com	Success	1:05 AM	1:05 AM	5m 60s
1:05 AM	ec2-54-242-188-185.compute-1.amazonaws.com	Success	1:05 AM	1:05 AM	5m 42s
1:05 AM	ec2-50-16-81-139.compute-1.amazonaws.com	Success	1:05 AM	1:05 AM	5m 53s
1:05 AM	ec2-23-21-46-96.compute-1.amazonaws.com	Success	1:05 AM	1:05 AM	5m 45s
1:05 AM	ec2-184-73-131-182.compute-1.amazonaws.com	Success	1:05 AM	1:05 AM	5m 56s
1:05 AM	ec2-54-242-249-112.compute-1.amazonaws.com	Success	1:05 AM	1:05 AM	5m 43s
1:05 AM	ec2-54-242-186-138.compute-1.amazonaws.com	Success	1:05 AM	1:05 AM	6m 2s
1:05 AM	ec2-72-44-44-107.compute-1.amazonaws.com	Success	1:05 AM	1:05 AM	5m 55s
1:05 AM	ec2-54-242-254-153.compute-1.amazonaws.com	Success	1:05 AM	1:05 AM	5m 32s
1:05 AM	ec2-184-73-18-17.compute-1.amazonaws.com	Success	1:05 AM	1:05 AM	5m 53s
1:05 AM	ec2-54-234-78-143.compute-1.amazonaws.com	Success	1:05 AM	1:05 AM	5m 47s
1:05 AM	ec2-50-16-187-188.compute-1.amazonaws.com	Success	1:05 AM	1:05 AM	6m 0s
1:05 AM	ec2-174-129-96-193.compute-1.amazonaws.com	Success	1:05 AM	1:05 AM	5m 57s
1:05 AM	ec2-54-242-255-224.compute-1.amazonaws.com	Success	1:05 AM	1:05 AM	5m 46s
1:05 AM	ec2-54-234-135-20.compute-1.amazonaws.com	Success	1:05 AM	1:05 AM	5m 27s
1:05 AM	ec2-72-44-42-32.compute-1.amazonaws.com	Success	1:05 AM	1:05 AM	5m 26s
1:05 AM	ec2-23-20-229-248.compute-1.amazonaws.com	Success	1:05 AM	1:05 AM	5m 51s
1:05 AM	ec2-184-73-211-202.compute-1.amazonaws.com	Success	1:05 AM	1:05 AM	5m 36s
1:05 AM	ec2-204-228-124-95.compute-1.amazonaws.com	Success	1:05 AM	1:05 AM	5m 23s
1:05 AM	ec2-54-234-78-9.compute-1.amazonaws.com	Success	1:05 AM	1:05 AM	5m 18s
1:05 AM	ec2-23-22-64-187.compute-1.amazonaws.com	Success	1:05 AM	1:05 AM	5m 48s
1:05 AM	ec2-23-22-241-155.compute-1.amazonaws.com	Success	1:05 AM	1:05 AM	5m 36s
1:05 AM	ec2-50-17-127-209.compute-1.amazonaws.com	Success	1:05 AM	1:05 AM	5m 16s
1:05 AM	ec2-54-242-252-185.compute-1.amazonaws.com	Success	1:05 AM	1:05 AM	5m 18s
1:05 AM	ec2-54-242-239-2.compute-1.amazonaws.com	Success	1:05 AM	1:05 AM	5m 38s
1:05 AM	ec2-23-20-29-123.compute-1.amazonaws.com	Success	1:05 AM	1:05 AM	5m 34s
1:05 AM	ec2-67-202-75-201.compute-1.amazonaws.com	Success	1:05 AM	1:05 AM	5m 32s
1:05 AM	ec2-72-44-40-171.compute-1.amazonaws.com	Success	1:05 AM	1:05 AM	5m 21s
1:05 AM	ec2-184-73-149-210.compute-1.amazonaws.com	Success	1:05 AM	1:05 AM	5m 29s
1:05 AM	ec2-54-242-212-25.compute-1.amazonaws.com	Success	1:05 AM	1:05 AM	5m 28s
1:05 AM	ec2-54-242-236-222.compute-1.amazonaws.com	Success	1:05 AM	1:05 AM	5m 29s
1:05 AM	ec2-54-242-236-222.compute-1.amazonaws.com	Success	1:05 AM	1:05 AM	5m 29s

Converge status by host



Design principle: Planning to handle failure is not planning to fail nor failing to plan

- Wrapper checks to see if its result is present and correct.
- There are a lot of moving parts. Different things break at different scales.
- Testing is essential but you'll always find new issues when running at scale.
- Data is stale.
- Make sure you have enough file handles!
- HTCondor can be overwhelmed by too many short jobs.
- Spots fail: HTCondor is designed to handle this.

Additional advice

- Keep tight with your friends. (Keep your friends close and your enemies closer.)
- DAGMan is your friend
 - Even when there aren't dependencies between jobs
- CycleServer is your friend
 - What the heck is going on?
 - The race: Jason wins.
- **Additional advice: maintain flexibility, balance**
 - Keep it simple
 - Throw stuff out
 - Elegant job wrapper with cached data
 - Keep it fun



Thank you, Questions?

- Utility Supercomputing 50 to 50,000 cores
- Visualization, Reporting
- Data scheduling: internal ↔ cloud
- Workload portability

