Putting Eggs in Many Baskets: Data Considerations in the Cloud

Rob Futrick, CTO



We believe utility access to technical computing power accelerates discovery & invention

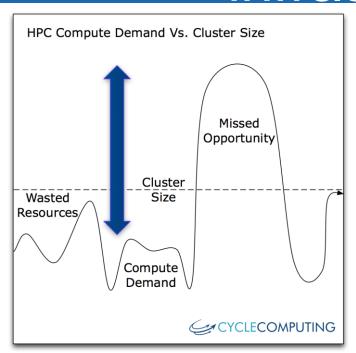


The Innovation Bottleneck:

Scientists/Engineers
forced to size their work to the
infrastructure their organization bought



Limitations of fixed infrastructure



Too small when needed most, Too large every other time...

- Upfront CapEx anchors you to aging servers
- Costly administration
- Miss opportunities to do better risk management, product design, science, engineering



Our mission: Write software to make utility technical computing easy for anyone, on any resources, at any scale



As an example...



Many users use 40 - 4000 cores, but let's talk about an example:

World's first PetaFLOPS
((Rmax+Rpeak)/2)
Throughput Cloud Cluster



Study Rules Out Global Warming Being a Natural Fluctuation With 99% Certainty

Posted by Soulskill on Saturday April 12, 2014 @11:39AM from the let's-blame-the-dinosaurs dept.



An anonymous reader writes

"A study out of McGill University sought to examine historical temperature data going back 500 years in order to determine the likelihood that global warming was caused by natural fluctuations in the earth's climate. The study concluded there was less than a 1% chance the warming could be attributed to simple fluctuations. The climate reconstructions take into account a variety of gauges found in nature, such as tree rings, ice cores, and lake sediments. And the fluctuation-analysis techniques make it possible to understand the temperature variations over wide ranges of time scales. For the industrial era, Lovejoy's analysis uses carbon-dioxide from the burning of fossil fuels as a proxy for all man-made climate influences – a simplification justified by the tight relationship between global economic activity and the emission of greenhouse gases and particulate pollution, he says. ... His study [also] predicts, with 95% confidence, that a doubling of carbon-dioxide levels in the atmosphere would cause the climate to warm by between 2.5 and 4.2 degrees Celsius. That range is more precise than – but in line with — the IPCC's prediction that temperatures would rise by 1.5 to 4.5 degrees Celsius if CO2 concentrations double."

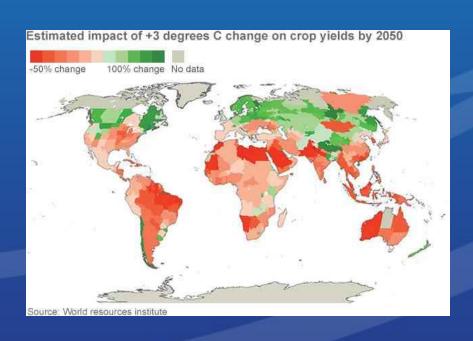
Read the 849 comments



x news x statistics x collapse x cuethedeniers x climatechange story



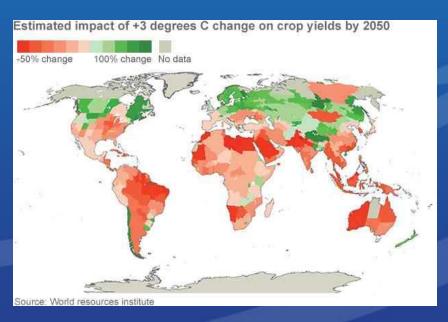
What do you think?



- Much of the worlds "bread basket" land will be hotter and drier
- Ocean warming is decreasing fish populations / catches



First, buy land in Canada?







Solar Energy

Wind power

Nuclear Fission energy

GeoThermal

Sure! But there have to be engineer-able solutions too.

Climate Engineering

Nuclear Fusion

BioFuels



Designing Solar Materials

The challenge is efficiency - turning photons to electricity

The number of possible materials is limitless:

- Need to separate the right compounds from the useless ones
- Researcher Mark Thompson, PhD:

"If the 20th century was the century of silicon, the 21st will be all organic. Problem is, how do we find the right material without spending the entire 21st century looking for it?"



Needle in a Haystack Challenge:

205,000 compounds totaling 2,312,959 core-hours or 264 core-years



205,000 molecules 264 years of computing

16,788 Spot Instances, 156,314 cores



205,000 molecules 264 years of computing

156,314 cores = 1.21 PetaFLOPS (~Rpeak)

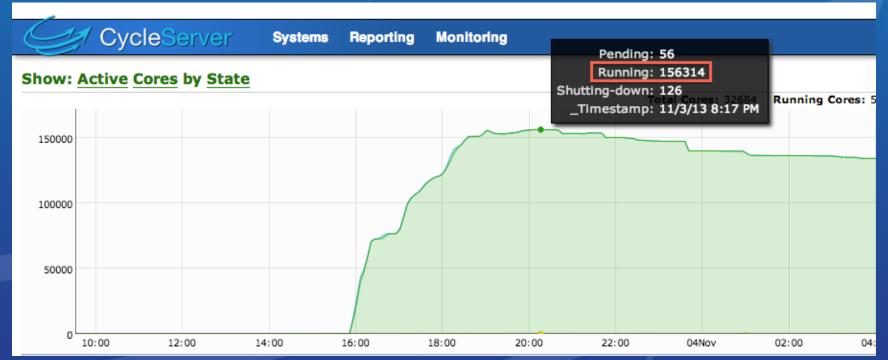


8-Region Deployment





1.21 PetaFLOPS (Rmax+Rpeak)/2, 156,314 cores

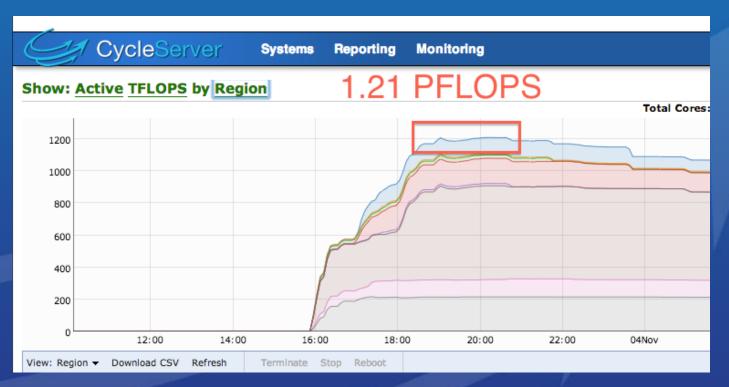




Each individual task was MPI, using a single, entire machine



Benchmark individual machines

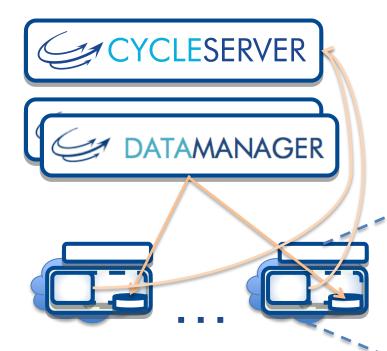




205,000 molecules 264 years of computing

Done in 18 hours Access to \$68M system for \$33k

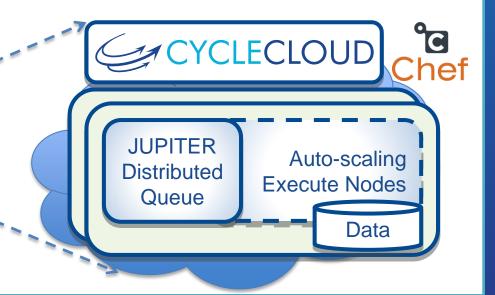




14 nodes controlling 16,788

How did we do this?

Automated in 8 Cloud Regions, 5 continents, double resiliency



Now Dr. Mark Thompson

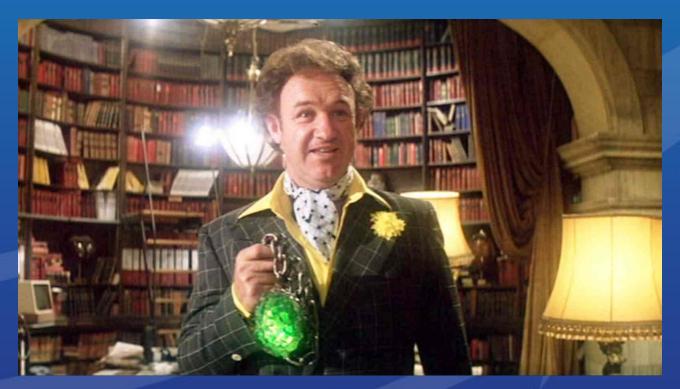


Is 264 compute
 years closer to
 making efficient solar
 a reality using organic
 semiconductors

Important?



Not to me anymore;)





Large Memory

Interconnect Sensitive

Whole sample set analysis

Large Grid MPI Runs

We see across all workloads

Needle in a Haystack runs

Interactive, SOA

Hi I/O Big Data runs



Users want to decrease Lean manufacturing's 'Cycle time'

(prep time + queue time + run time)



Everyone we work with faces this problem





External resources (Open Science Grid, cloud) offer a solution!





How do we get there?



In particular how do we deal with data in our workflows?



Several of the options...

- Local disk
- Object Store / Simple Storage Service (S3)
- Shared FS
 - NFS
 - Parallel file system (Lustre, Gluster)
- NoSQL DB



Let's do this through examples



Compendia BioSciences (Local Disk)





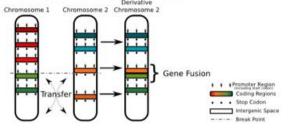
Cloud Computing Enabled Discovery of Gene Fusions from Whole Transcriptome NGS Tumor Samples

The Cancer Genome Atlas



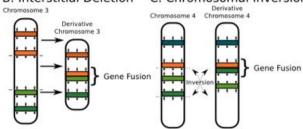
Gene Fusion Overview

A. Chromosomal Translocation



B. Interstitial Deletion

C. Chromosomal Inversion



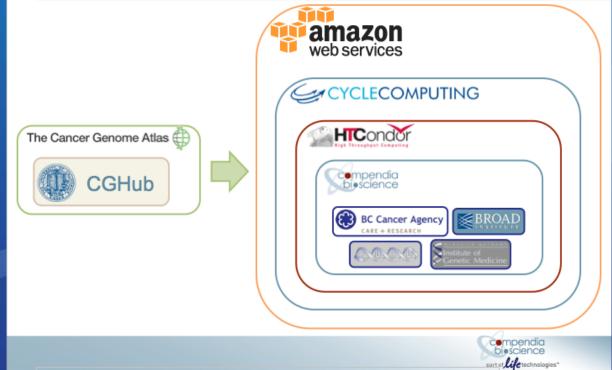
http://en.wikipedia.org/wiki/Gene_fusion

Click to add text





Fusion Calling Platform

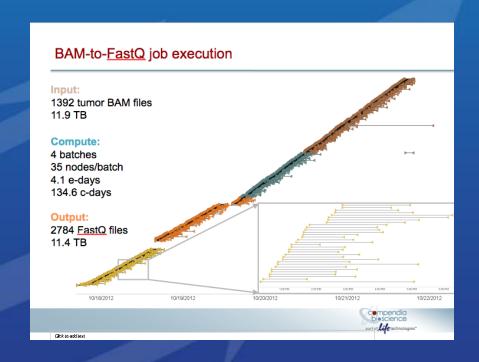






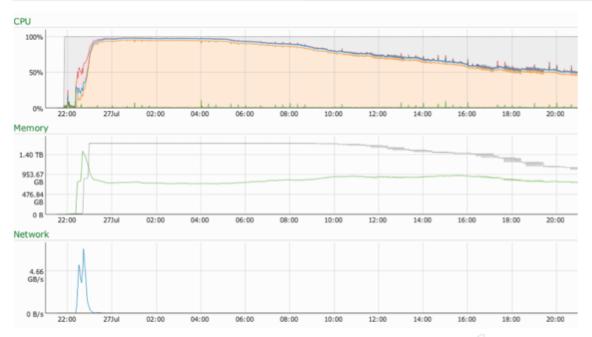
As described at TriMed Conference

- Stream data out of TCGA into S3 & Local machines (concurrently on 8k cores)
- Run analysis and then place results in S3
- No Shared FS, but all nodes are up for a long time downloading (8000 cores * xfer)





Defuse node utilization

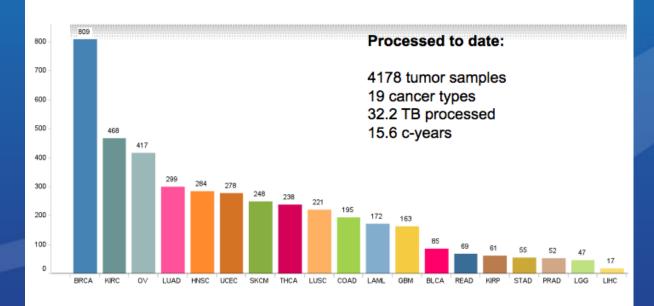




Click to add text



Fusion Caller Processing Summary



As of 12/31/2012

Click to add text





Pros & Cons of Local Disk

Pros:

- Typically no application changes
- Data encryption is easy; No shared key management
- Highest speed access to data once local

Cons:

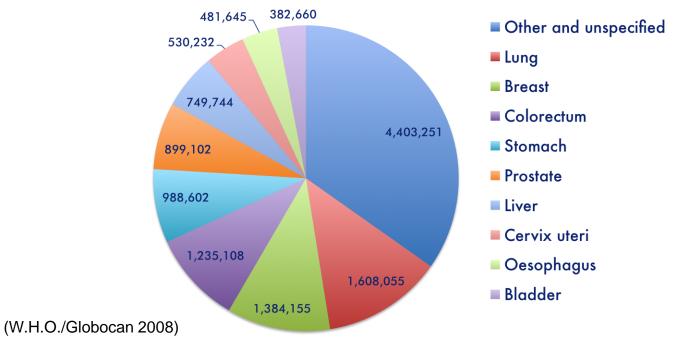
- Performance cost/impact of all nodes downloading data
- Only works for completely mutually exclusive workflows
- Possibly more expensive;
 Transferring to a shared filer and then running may be cheaper



Novartis (S3 sandboxes)



New Cases of Cancer per Year: 12.66 Million

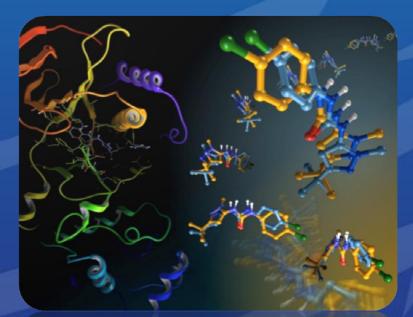




Every day is crucial and costly



Challenge: Novartis run 341,700 hours of docking against a cancer target; impossible to do in-house.





Most Recent Utility Supercomputer server count:

```
$ knife node list | wc -l
    10599
$ ■
```



AWS Console view (the only part that rendered):

Request Spot Instances	Cancel	Pricing History		C	•	8)
Viewing: Active + Se	arch		I< <	1 to 50 of 10598 Ite	ems	>	> I



Cycle Computing's cluster view:

	ycle		Systems	Reporting	Monitoring					
nef overvi	iew for	chef-serve	r-11.cyccld	.com						
Current host	stats				Converge stats (las	st hour)	Alerts			
# Chef Servers: 1 Total Converges: 3944										
	# Hosts:				Successful Converg		♣ Fr Feb 01 2013 18:21:53 GHT-0500 (EST): ec2-23-22-131-239.compute-1.amazonaws.com failed to converge			
# Converged Hosts: 10312 Failed Converges: 92						es: 92				
# Unconverged Hosts: 31							Converge status by host			
tecent conve	erges									
Download CSV			Mark Persistent			Search:				
	A Host Nan			Start Time	End Time	Duration				
No. of Concession, Name of Street, or other Persons, Name of Street, or ot		4-122-169.comp -125-129.compu		Parket Service	P. W. Ser	12m 6s 6m 19s				
The Per	902-54-24	2-78-85.compute	Success	how the	Trans. West	6m 40s				
	952-54-24	2-93-226.comput	Success	Color Service	Trade Ann	6m 22s				
yes deer		3-138-253.comp		high life	From Bet	6m 36s				
to the		2-240-184.comps		Trape Mile	11.00 PM	6m 17s 6m 20s				
on the		2-144-205.comp -234-71.compute		Train Per	71.00 PM	6m 20s 6m 9s				
the later		2-83-75.compute		Total Service	100.00	6m 1s				
the deep		-29-217.compute		From Sec.	h par det	6m 4s				
de ter		-169-52.compute		TUDE DEL	Track Ber	6m 13s 6m 20s				
on the		2-44-228 comput 2-59-145 comput		Train Mer	11.00 PM	6m 4s				
and their		2-202-232.comp		DOM: NO	Proper work	6m 20s				
jes der	ec2-67-20	2-21-38 compute	Success	NUMBER	Trust, Med	6m 4s				
Op. me		29-61-233.comp		Tubs Hert	(-, g) eer	5m 55s				
30 PM		29-137-142.com 3-14-15.compute		high Mr	hope det	5m 47s 6m 10s				
No. see		-137-89.compute		7-25 PM	7.00 800	6m 2s				
rus Pre		2-240-75.comput	Success	Number	From Per	5m 60s				
-		2-188-185.comp		1.05.00	Supplement	5m 42s				
Jos. Ber.		-41-139 compute -6-96 compute-1		TOTAL SERVICE	11.00 000	5m 53s 5m 45s				
-0.00		2-181-52.comput		149.00	State Ben	5m 56s				
The Pier		2-249-112.comp		high Phi	10,000 894	5m 43s				
in an		2-186-138.comp		7-26 886	Print, men	6m 2s				
35 Per		2-254-153 compute		Transport	11.00 PM	5m 55s 5m 32s				
St. Are		3-18-17.compute		Trube Arm	has been	5m 53s				
top the		4-78-143.comput	Success	hops the	house were	5m 47s				
in sec		-187-188.comput		YOR ME	71 M MW	6m 0s				
District		29-96-193 comps 2-255-224 comps		Color Medi	TOTAL PROPERTY.	5m 57s 5m 46s				
To Are		4-135-20 comput		Total Ber	1.00 80	5m 46s 5m 27s				
70 MT	ec2-72-44	-42-32.compute-	Success	Frage AM	7-36-99	5m 26s				
3.00		-229-248.comput		1.25 800	From Mee	5m 51s				
to the		2-211-202 comp 36-194-94 comp		has the	71.00 PM	5m 36s 5m 23s				
St. eer		4-28-9 compute-		Color Per	1, 30 MH	5m 23s 5m 18s				
No wee	ec2-23-22	-64-147.compute	Success	Fraga Area	Project war.	5m 48s				
de dec		-241-155.comput		Tube deep	PLOS BM	5m 36s				
State Per		-127-209.comput		Total Per	11,00,000	5m 16s 5m 18s				
0.00		2-252-185.comp 2-239-2.compute		CUE DE	has the	5m 18s 5m 38s				
10.00		-29-125,compute	Success	Trup Met	Total Ber	5m 34s				
(in the	ec2-67-20	2-7-201 compute	Success	Inches place	Print, Ber	5m 32s				
No. Per		-40-171.compute		100	11,000,000	5m 21s				
in the		3-149-210.comp- 2-212-25.comput		100	200	5m 29s 5m 28s				
On the		2-236-222.comp		7-26 MH	Print Net	5m 29s				
es Per	ec2-54-24	2-216-136.comps	Success	high plen	Transport	5m 21s				
20 000		2-103-28.comput		1-25 800	Track Person	5m 39s				
I'm Are	er2-50-19	-148-13 romoute	Sunnes	complex	2110 844	5m 164				



Metric	Count
Compute Hours of Science	341,700 hours
Compute Days of Science	14,238 days
Compute Years of Science	39 years
AWS Instance Count	10,600 instances

CycleCloud, HTCondor, & Cloud finished impossible workload, \$44 Million in servers, ...





39 years of drug design in 11 hours on 10,600 servers for < \$4,372



Does this lead to new drugs?



Novartis announced 3 new compounds going into screening based upon this 1 run



Pros & Cons of S3

Pros:

- Parallel scalability; high throughput access to data
- Only bottlenecks occur on S3 access
- Can easily and greatly increase capacity

Cons:

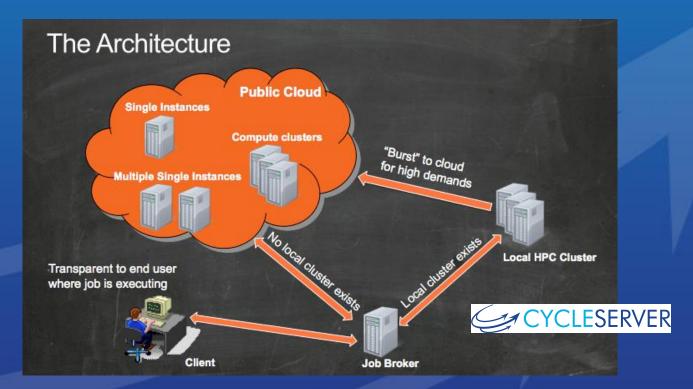
- Only works for completely mutually exclusive workflows
- Native S3 access requires application changes
- Non-native S3 access can be unstable
- Latency can affect performance



Johnson & Johnson (Single node NFS)

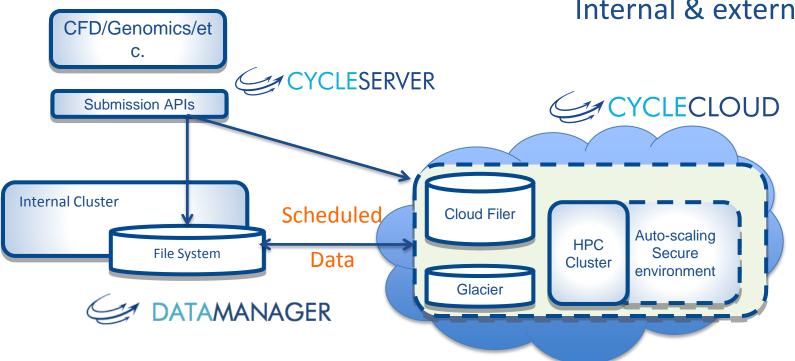


JnJ @ AWS re:Invent 2012





JnJ Burst use case Internal & external





COMPUTING

(Patents Pending)

Pros & Cons of NFS

Pros:

- Typically no application changes
- Cheapest at small scale
- Easy to encrypt data
- Performance great at (small) scale and/or under some access patterns
- Great platform support

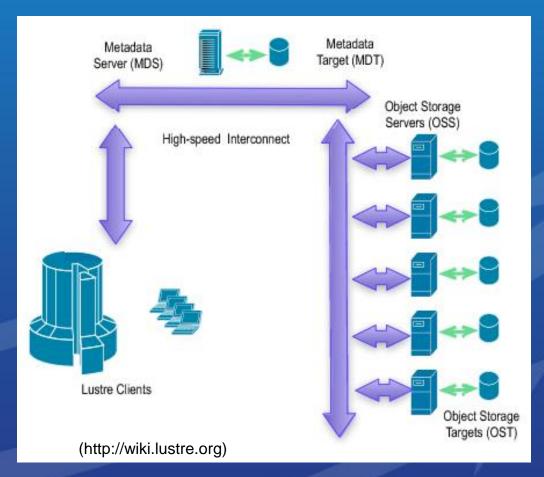
Cons:

- Filer can easily become performance bottleneck
- Not easily expandable
- Not fault tolerant; Single point of failure
- Backup and recovery



Parallel Filesystem (Gluster, Lustre)





Pros & Cons of Parallel FS

Pros:

- Easily expand capacity
- Read performance scalability
- Data integrity

Cons:

- Greatly increased administration complexity
- Performance for "small" files can be atrocious
- Poor platform support
- Data integrity and backup
- Still has single points of failure



NoSQL DB (Cassandra, MongoDB)















(http://strata.oreilly.com/2012/02/nosql-non-relational-database.html)



Pros & Cons of NoSQL DB

Pros:

- Best performance for appropriate data sets
- Data backup and integrity
- Good platform support

Cons:

- Only appropriate for certain data sets and access patterns
- Requires application changes
- Application developer and Administration complexity
- Potential single point of failure



Large Memory

Interconnect Sensitive

Whole sample set analysis

Large Grid MPI Runs

That's a survey of different workloads

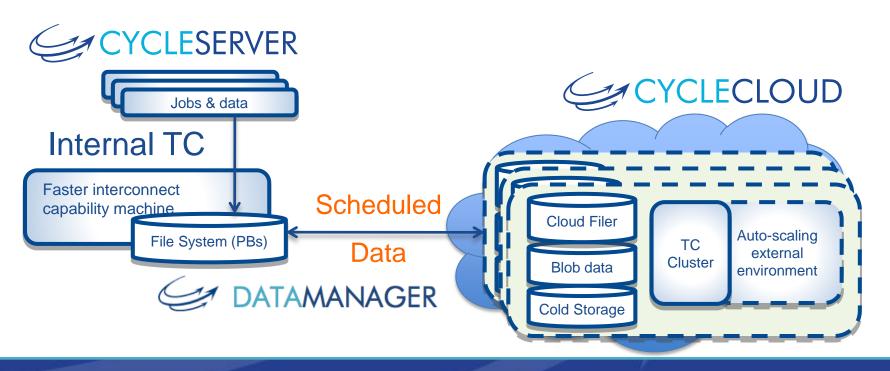
Needle in a Haystack runs

Interactive, SOA

Hi I/O Big Data runs



Depending upon your use case



There are a lot more examples...



Life Sciences

- 39.5 years of drug compound computations in 9 hours, at a total cost of \$4,372
- 10,000 server cluster seamlessly spanned US/EU regions
- Advanced 3 new otherwise unknown compounds in wet lab

SCHRÖDINGER.

Life Sciences

- 156,000-core utility supercomputer in the Cloud
- Used \$68M in servers for 18 hours for \$33,000
- Simulated 205,000 materials (264 years) in 18 hours

Nuclear engineering Utilities / Energy

- Approximately 600-core MPI workloads run in Cloud
- Ran workloads in months rather than years
- Introduction to production in 3 weeks



Rocket Design & Simulation

- Moving HPC workloads to cloud for burst capacity
- Amazon GovCloud, Security, and key management
- Up to 1000s of cores for burst / agility



Asset & Liability Modeling (Insurance / Risk Mgmt)

- Completes monthly/quarterly runs 5-10x faster
- Use 1600 cores in the cloud to shorten elapsed run time from ~10 days to ~ 1-2 days

Johnson Johnson Manufacturing & Design

- Enable end user on-demand access to 1000s of cores
- Avoid cost of buying new servers
- Accelerated science and CAD/CAM process



Hopefully you see various data placement options



Each have pros and cons...

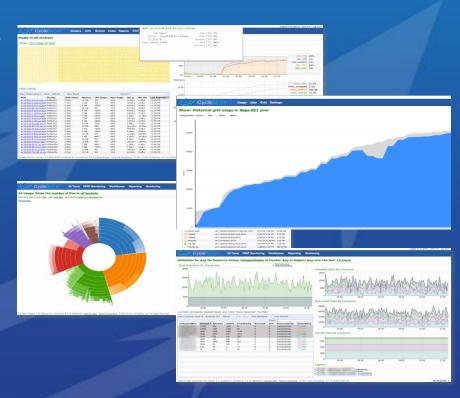
- Local disk
- Simple Storage Service (S3)
- Shared FS
 - NFS
 - Parallel file system (Lustre, Gluster)
- NoSQL DB



We write software to do this...

Cycle easily orchestrates workloads and data access to local and Cloud TC

- Scales from 100 100,000's of cores
- Handles errors, reliability
- Schedules data movement
- Secures, encrypts and audits
- Provides reporting and chargeback
- Automates spot bidding
- Supports Enterprise operations





Does this resonate with you?



We're hiring like crazy: Software developers, HPC engineers, devops, sales, etc.

jobs@ cyclecomputing.com



Now hopefully...



You'll keep these tools in mind



as you use HTCondor



to accelerate your science!

