

IceProd Scalability

Software & Computing Science Advisory Panel
January 27, 2021



Overview

- IceProd Multi-User Support
- Workflow Management
- Cloudburst Runs
- IceCube Upgrade & Gen2 Support

IceProd Multi-User Support

2018-4

Multi-user available for ~2 years

- Several WGs using IceProd for official processing
 - Usually the tech lead or other designated submitter
- Official simulations more diversified as a result
 - WGs run their own requested simulations, with coordination from SimProd group

Future of multi-user support

- A simpler interface for the average analyser
 - Designed for a simple script with an input and output file
- Better error messages
 - Hide internal technical details and put in plain language
 - Guide towards solution for common problems

Workflow Management

2018-9

Scheduling

- Removed one wrapper pilot due to insufficient maintenance effort
 - Less control over placement, but seems more robust
- Improvements for supercomputer environments, or network-restricted sites
 - Allow out-of-job data handling

Monitoring

- Improvements in tracking where and what is running
 - Grafana dashboards tracking usage by site, dataset, job type

Common theme: lack of effort resulted in slow progress

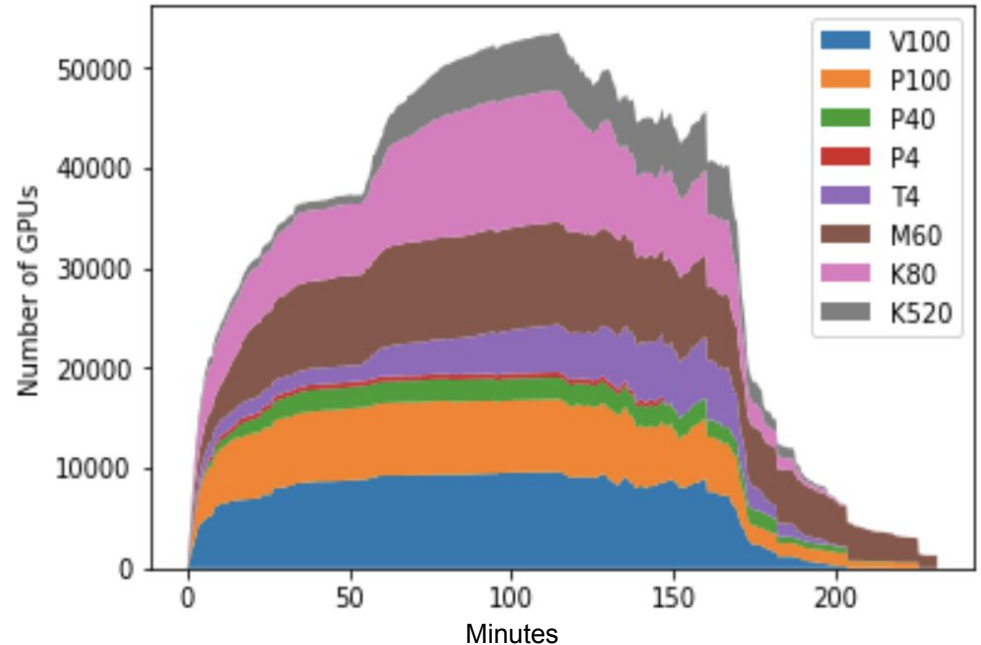
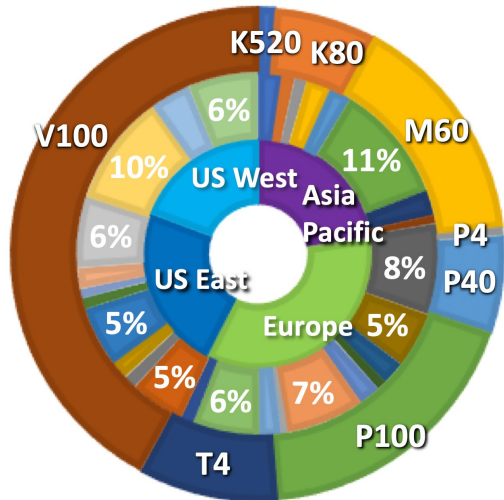
- Lots of work still to do in these areas

Cloudburst Runs

3 separate runs to test pre-exascale compute in the cloud

2018-9
2018-10

- 1 run for peak FLOPs:
All the GPUs we could buy,
compute-intensive workload



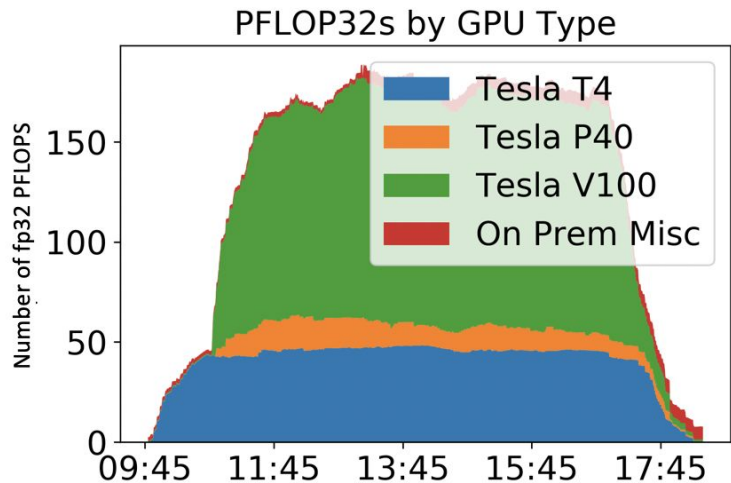
Cloudburst Runs

3 separate runs to test pre-exascale compute in the cloud

2018-9
2018-10

- 1 run for peak FLOPs:
All the GPUs we could buy,
compute-intensive workload
- 1 “economical” run:
Only use spot instances for 3 most
efficient GPU types

	PFLOP32h	Jobs	Cost (+- 15%)
All	1082	151k	
Cloud	1033	145k	\$60k
T4	316	48k	\$9k
T4 Fraction	30%	33%	15%

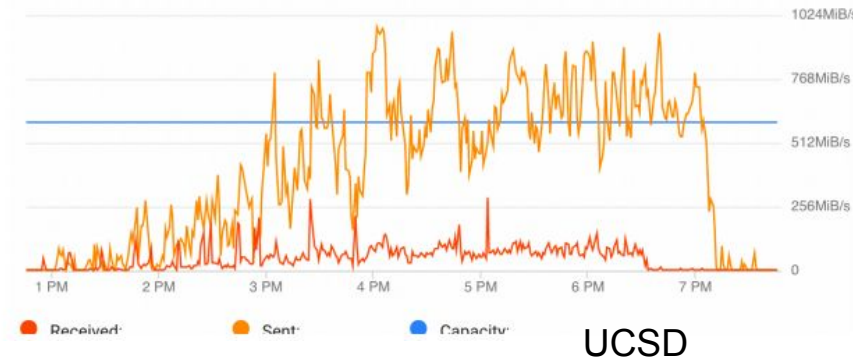
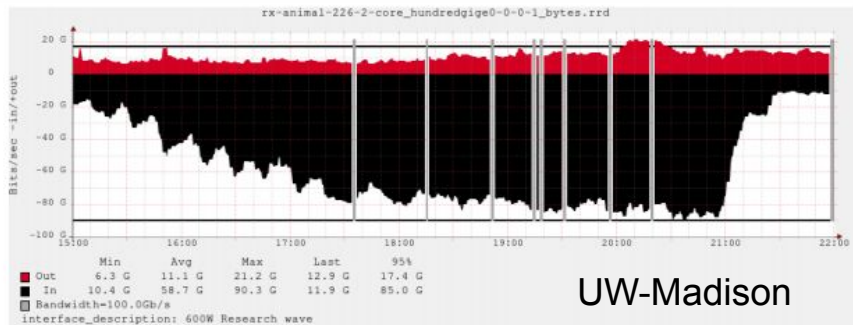


Cloudburst Runs

3 separate runs to test pre-exascale compute in the cloud

2018-9
2018-10

- 1 run for peak FLOPs:
All the GPUs we could buy,
compute-intensive workload
- 1 “economical” run:
Only use spot instances for 3 most
efficient GPU types
- 1 “data intense” run
Large outfile workflow to test
networking



Cloudburst Runs

A good scaling exercise

- >50k running jobs, >1M idle jobs in HTCondor queue
- Stress test of IceProd
 - Fixed a number of scaling bottlenecks
 - Improved handling of varied network connectivity
- Storage can do 100 Gbps writes
- Able to run this in parallel with normal production

2018-9
2018-10

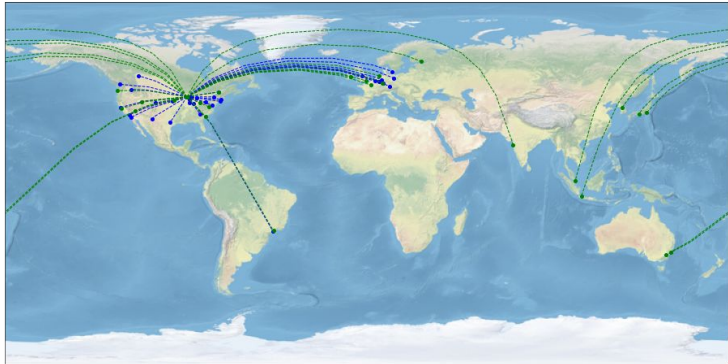
Demonstrated support for >10x resource usage

IceCube Upgrade & Gen2 Support

2021-3

Predictions

- Simulation needs grow by a large factor
 - Increased volume, more sensors
- Data reprocessing needs increase
- Analysis computing continues to grow



Scaling our Compute

- More allocations / funding to grow compute pool
- Better monitoring and scheduling to reduce waste
 - Site/node failure handling
 - Steering jobs to optimal sites
- Start to look at data movement, asynchronous i/o

Backup Slides

What is IceProd?

Data provenance

- Configuration for how a file was generated or processed
- Which software, what versions, when/where it ran, ...

Dataset submission

- Monitor job status, resource usage
- Retry failed jobs - resubmit with different requirements

Use cases:

- Simulation production
- Experiment data processing
- Common analysis processing
- Other large-scale workloads

Pyglidein

A python server-client pair for submitting HTCondor glidein jobs on remote batch systems.

<https://github.com/WIPACrepo/pyglidein>

Motivation / requirements:

- MFA
- Lightweight library for easy remote operation
- UNIX philosophy

