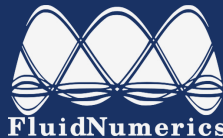


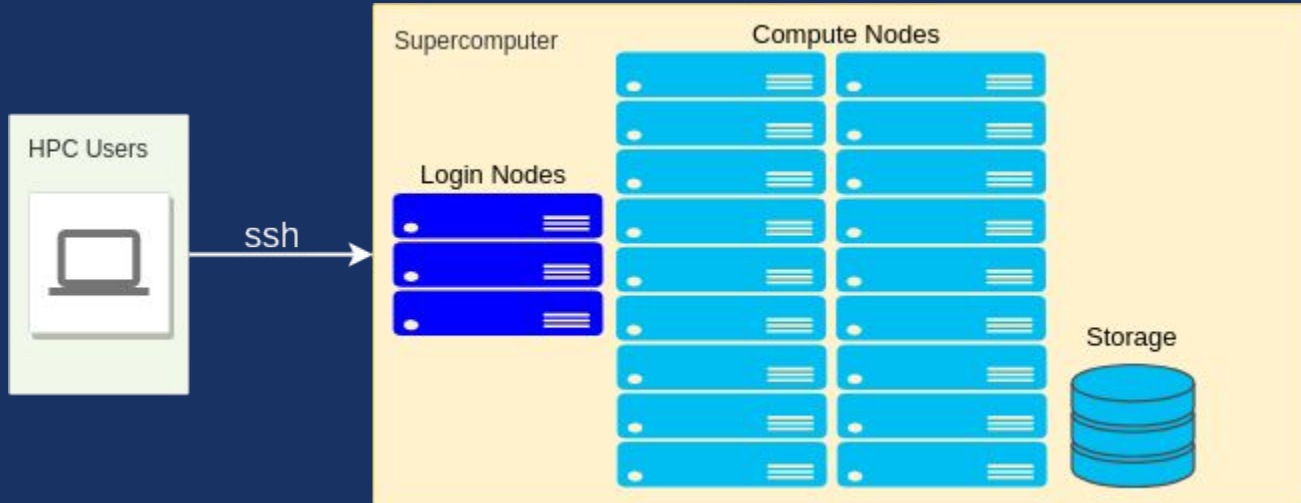
Bare-Metal Style HPC Clusters On Google Cloud Platform

Dr. Joseph Schoonover

*CEO, Cloud-HPC Systems Architect
Fluid Numerics, LLC
Boulder, CO*

joe@fluidnumerics.com



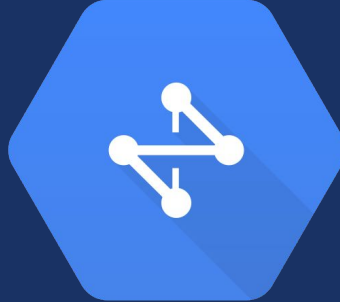


This talk will show how we can replicate this type of environment on the cloud

HPC Components



Compute



Network



Storage



Deployment Manager

<https://cloud.google.com/deployment-manager/>

Infrastructure as code

→ Python or Jinja templates
and YAML dictionaries

Resources:

```
- type: compute.v1.instance  
  name: quickstart-deployment-vm
```

Properties:

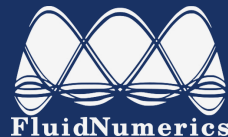
```
  zone: us-central1-f  
  machineType: <MACHINE TYPE>
```

Disks:

```
- deviceName: boot  
  type: PERSISTENT  
  boot: true  
  autoDelete: true  
  initializeParams:  
    sourceImage: <OS IMAGE>
```

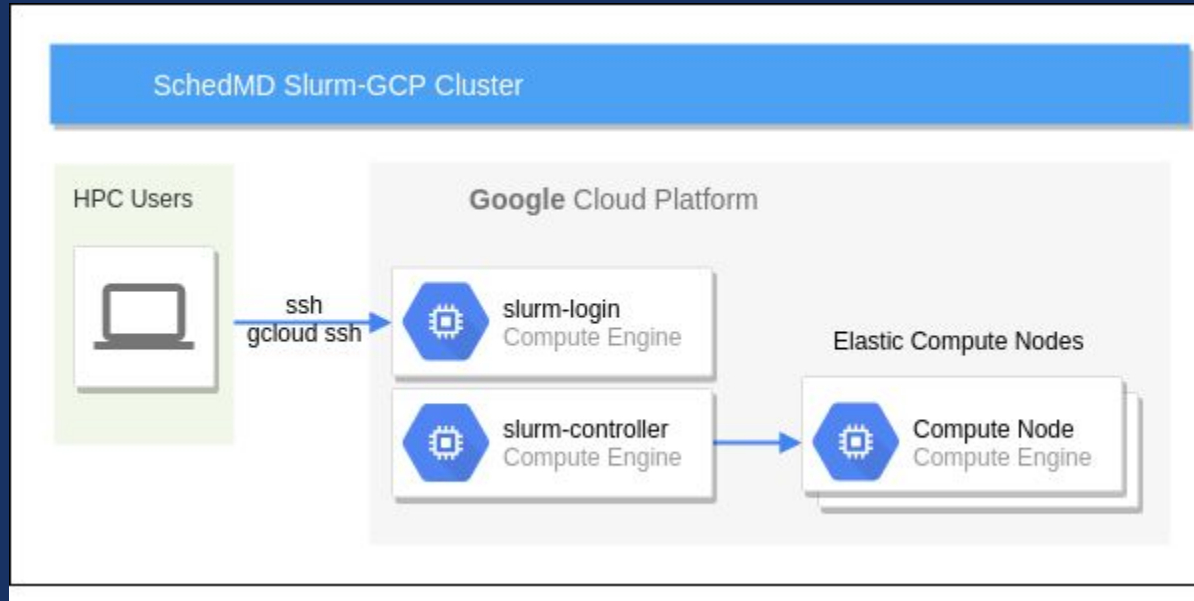
networkInterfaces:

```
- network: <NETWORK>
```



SchedMD : Slurm-GCP

<https://github.com/schedmd/slurm-gcp>

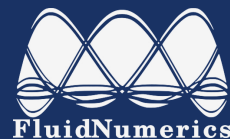


SchedMD : Slurm-GCP

How to use : Customize YAML file

```
#slurm-cluster.yaml
```

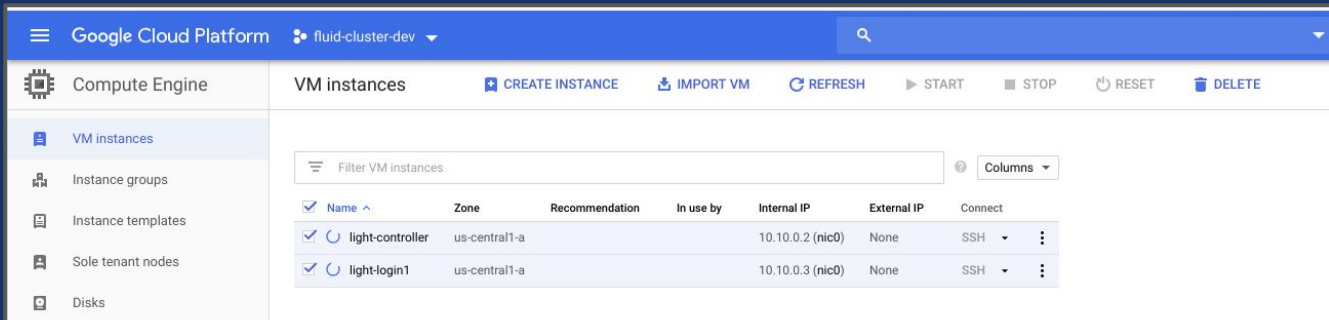
```
imports:
- path: slurm.jinja
resources:
- name: slurm-cluster
  type: slurm.jinja
properties:
  cluster_name           : g1
  static_node_count      : 2
  max_node_count         : 10
  zone                   : us-central1-b
  region                 : us-central1
  cidr                   : 10.10.0.0/16
  controller_machine_type : n1-standard-2
  compute_machine_type   : n1-standard-2
  login_machine_type     : n1-standard-2
```



SchedMD : Slurm-GCP

How to use : Deploy

```
gcloud deployment-manager deployment create slurm-cluster --project=<project id> --config=slurm-cluster.yaml
```



The screenshot displays the Google Cloud Platform console interface for the 'fluid-cluster-dev' project. The left sidebar shows the 'Compute Engine' section with a list of resources: 'VM instances' (selected), 'Instance groups', 'Instance templates', 'Sole tenant nodes', and 'Disks'. The main content area is titled 'VM instances' and includes action buttons: 'CREATE INSTANCE', 'IMPORT VM', 'REFRESH', 'START', 'STOP', 'RESET', and 'DELETE'. Below these buttons is a search bar labeled 'Filter VM instances' and a 'Columns' dropdown menu. A table lists two VM instances:

<input checked="" type="checkbox"/>	Name ^	Zone	Recommendation	In use by	Internal IP	External IP	Connect
<input checked="" type="checkbox"/>	light-controller	us-central1-a			10.10.0.2 (nic0)	None	SSH ▾ ⋮
<input checked="" type="checkbox"/>	light-login1	us-central1-a			10.10.0.3 (nic0)	None	SSH ▾ ⋮

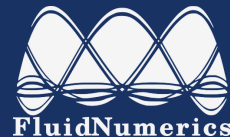


How to use : Login

```
ssh <user>@<ip-address>
```

[illegible][illegible]

```
/usr/bin/id: cannot find name for group ID 2078467674
[joeschoonover@light-login1 ~]$
```



Customizations

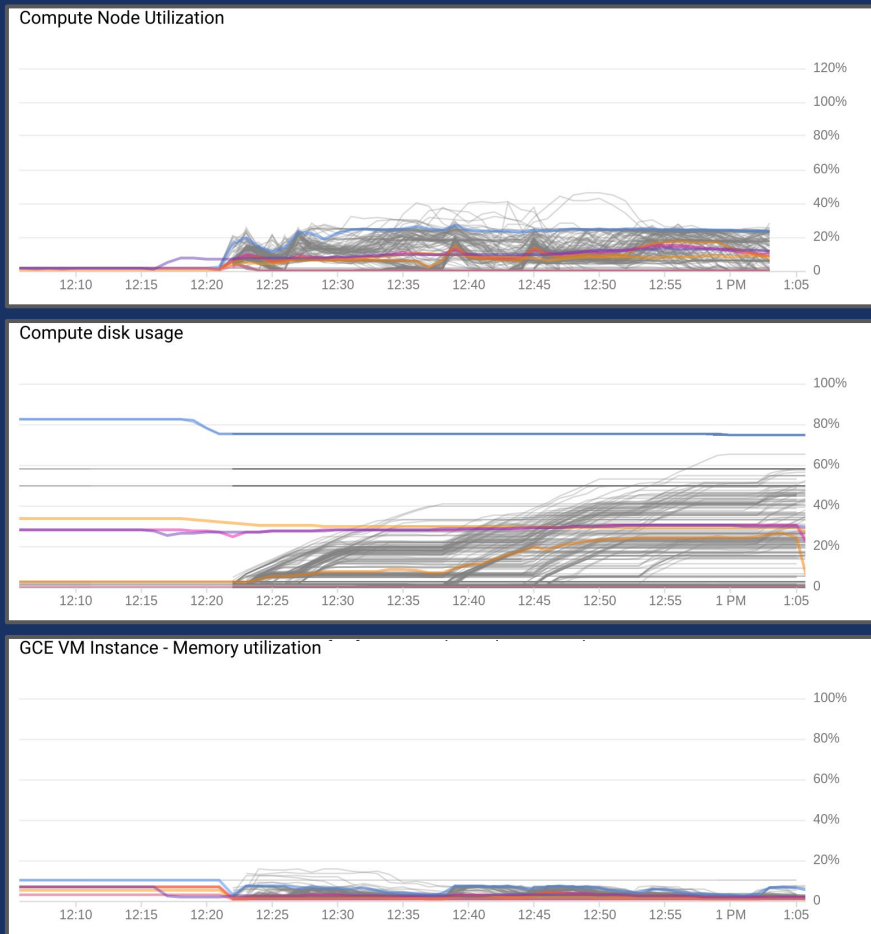
- Stackdriver monitoring integrations
- Multiple, user-locked, micro-login nodes
- Additional NFS storage and Lustre Storage
- Multiple partitions for multiple application/developer support
- Federation - burst expand on prem resources
- Python/Jupyter notebook integrations

Stackdriver

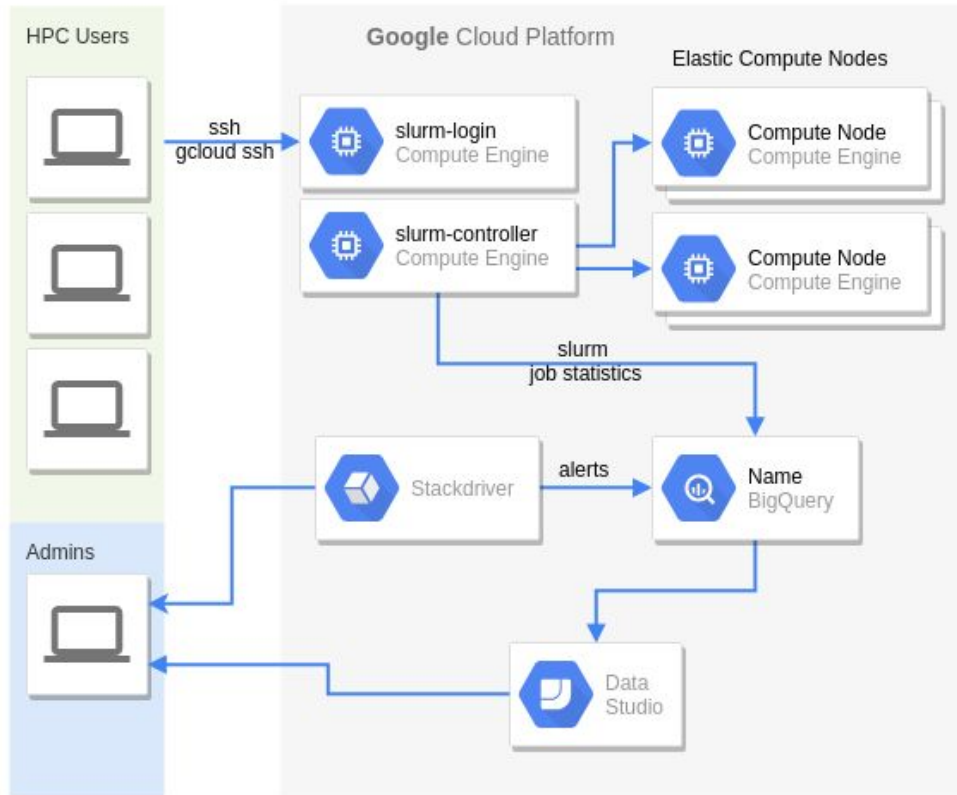
Built in system monitoring and alerting platform

Custom metrics can be integrated (e.g. GPU utilization)

This data can be used to resize compute nodes to maximize resource utilization and reduce costs



Architecture: Slurm Cluster



Compute Resources

Slurm Cluster

github.com/schedmd/slurm-gcp

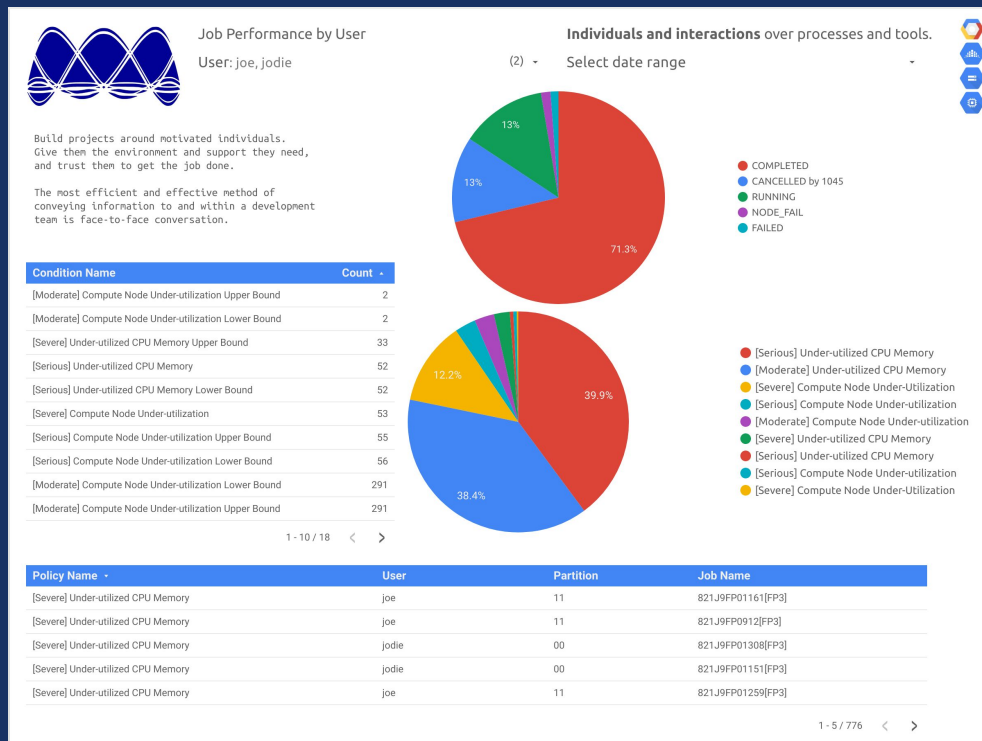
Customizations

+multi-partition +multi-zone

+environment modules

+github.com/spack/spack

Slurm + Stackdriver + BigQuery + DataStudio



Use Case

GPU Platform Comparisons

Partitions

highmem-64 (64 CPU + 416 GB RAM) + 8 Nvidia®
Tesla® V100 GPUs

highmem-32 (32 CPU + 208 GB RAM) + 4 Nvidia®
Tesla® P100 GPUs

highmem-16 (16 CPU + 104 GB RAM) + 4 Nvidia®
Tesla® P4 GPUs

highmem-16 (16 CPU + 104 GB RAM) + 4 Nvidia®
Tesla® K80 GPUs

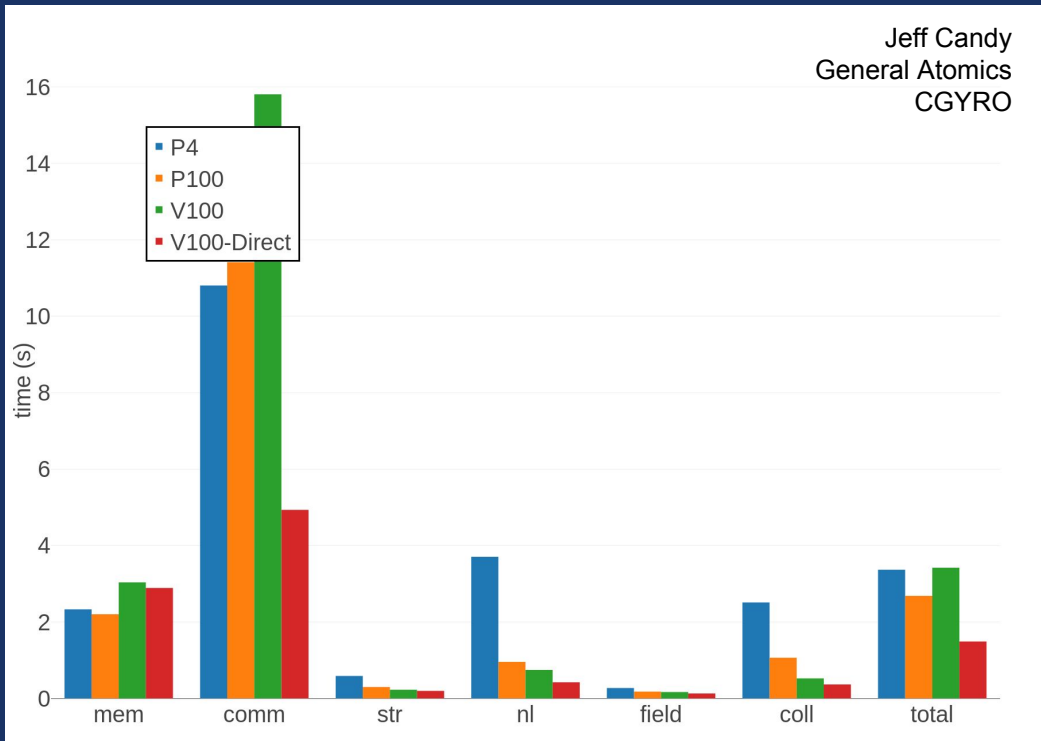
highmem-8 (8 CPU + 52 GB RAM) + 1 Nvidia®
Tesla® K80 GPU

highmem-8 (8 CPU + 52 GB RAM) + 1 Nvidia®
Tesla® P100 GPU

standard-32 (32 CPU + 120 GB RAM)

Use Case

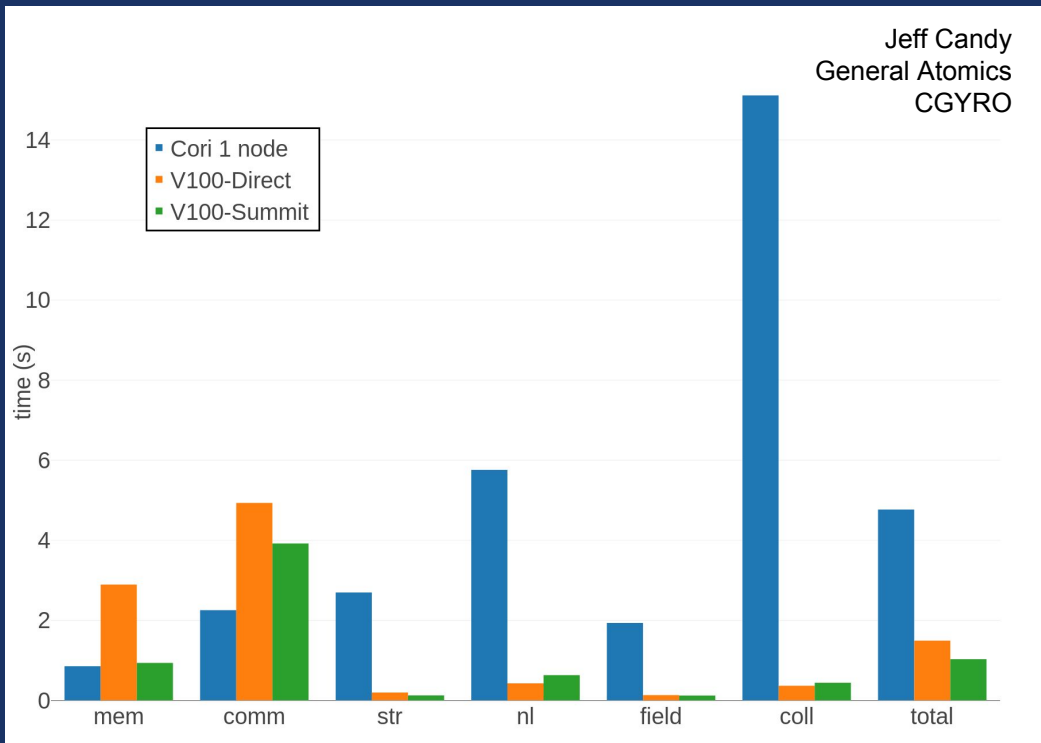
GPU Platform Comparisons



- Highmem-32 + 4xV100 GPU
- Highmem-32 + 4xP100 GPU
- Highmem-32 + 4xV100 GPU
- Highmem-32 + 4xV100 GPU + GPU-Direct MPI

Use Case

Cori, Summit, and GCP Comparisons



- Highmem-32 + 4xV100 GPU
- Highmem-32 + 4xP100 GPU
- Highmem-32 + 4xV100 GPU
- Highmem-32 + 4xV100 GPU + GPU-Direct MPI

Use Case

Cloud Cost Analysis

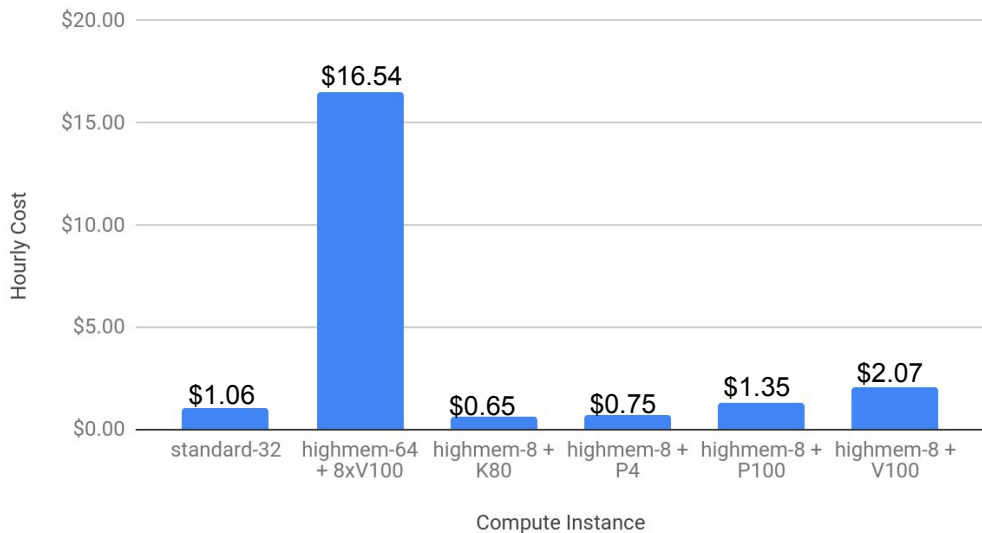


- Standard-32
- Highmem-8 + 1xP4 GPU
- Highmem-8 + 1xP100 GPU
- Highmem-8 + 1xV100 GPU
- Highmem-64 + 8xV100 GPU

Use Case

Cloud Cost Analysis

Hourly Cost vs. Compute Instance

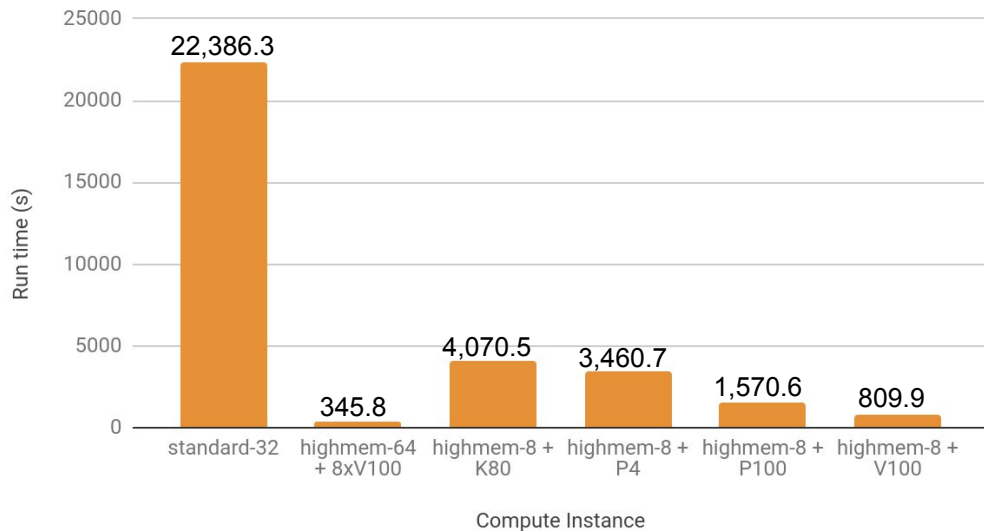


- Standard-32
- Highmem-8 + 1xP4 GPU
- Highmem-8 + 1xP100 GPU
- Highmem-8 + 1xV100 GPU
- Highmem-64 + 8xV100 GPU

Use Case

Cloud Cost Analysis

Run time (s) vs. Compute Instance

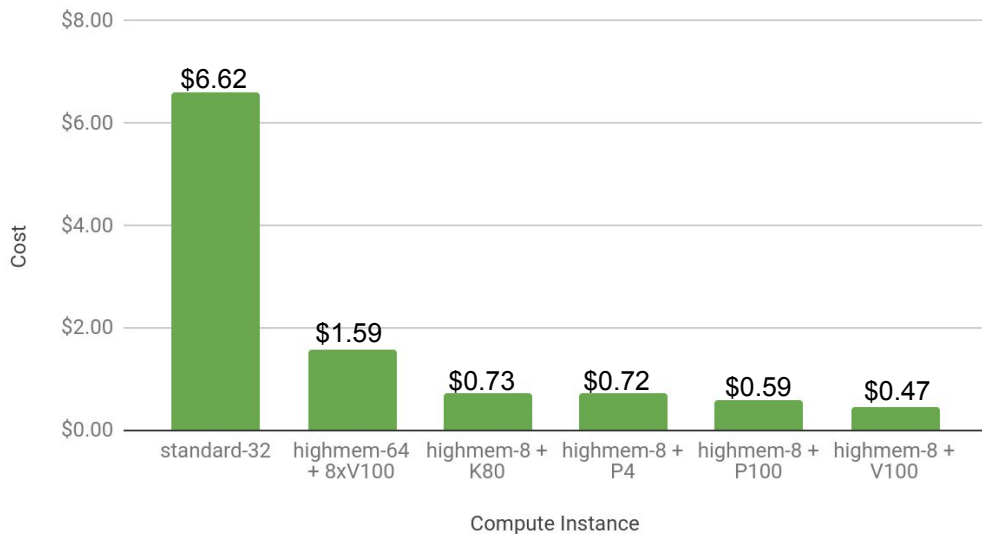


- Standard-32
- Highmem-8 + 1xP4 GPU
- Highmem-8 + 1xP100 GPU
- Highmem-8 + 1xV100 GPU
- Highmem-64 + 8xV100 GPU

Use Case

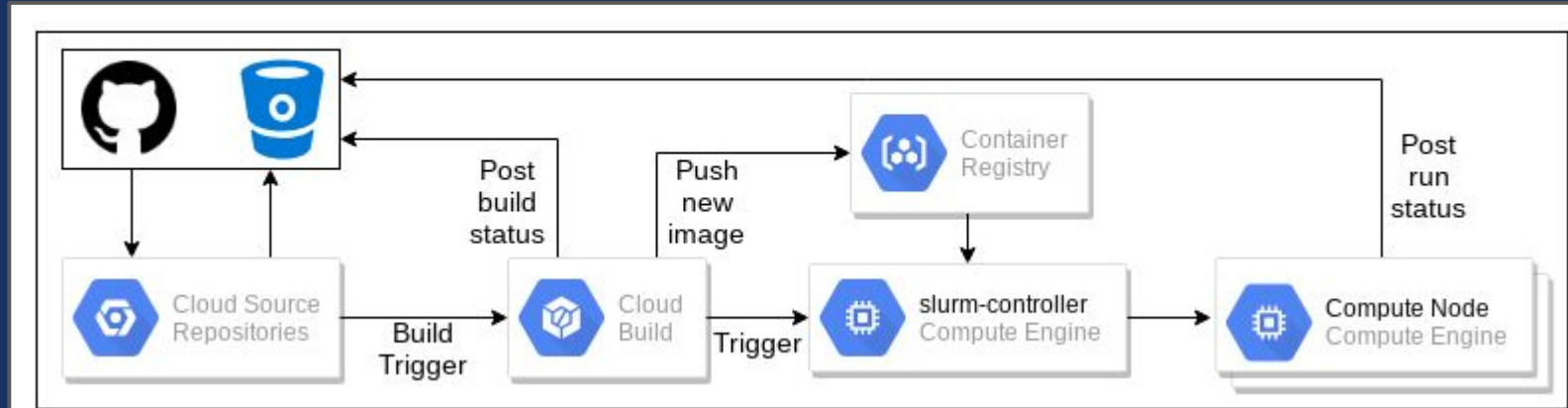
Cloud Cost Analysis

Cost vs. Compute Instance



- Standard-32
- Highmem-8 + 1xP4 GPU
- Highmem-8 + 1xP100 GPU
- Highmem-8 + 1xV100 GPU
- Highmem-64 + 8xV100 GPU

HPC DevOps



A Cultural Shift

Currently, organizations maintain one cluster that supports many teams

Teams of scientists/developers are separated from the infrastructure and admin teams

Developer teams or divisions can experiment and develop on their own cluster.

This would require infrastructure and admin individuals to interact more closely with scientists and developers

Impromptu Tutorial this afternoon

At 3:45 , we'll do a tutorial where you can spin up your own elastic slurm cluster on GCP.

<https://codelabs.developers.google.com/codelabs/hpc-slurm-on-gcp/>

Further Questions

Contact : joe@fluidnumerics.com

<https://fluidnumerics.com>