# Slurm in the Clouds

Nick Ihli
SchedMD

# Slurm User Group Meeting 2021

# Agenda

All times are US Mountain Daylight (UTC-6)

| Time | Speaker | Title |
|---|---|---|
| 9:00 - 9:50 | Jason Booth | Field Notes 5: From The Frontlines of Slurm Support |
| 10:00 - 10:25 | Nate Rini | REST API *and also* Containers |
| 10:30 - 10:50 | Marshall Garey | burst_buffer/lua and slurmscriptd |
| 11:00 - 11:25 | Nick Ihli | Slurm in the Clouds |
| 11:30 - 11:50 | Tim Wickberg | Slurm 21.08 and Beyond |

# Welcome

- Five separate presentations, five separate streams
- Presentations will remain available for at least two weeks after SLUG'21 concludes
- Presentations are available through the SchedMD Slurm YouTube channel
  - https://youtube.com/c/schedmdslurm
- Or through direct links from the agenda
  - https://slurm.schedmd.com/slurm_ug_agenda.html

# Asking questions

- Feel free to ask questions throughout through YouTube's chat
- Chat is moderated by SchedMD staff
  - Tim McMullan, Ben Roberts, and Tim Wickberg
  - Also identified by the little wrench symbol next to their name
- Questions will be relayed to the presenter by the moderators
  - Some may be deferred to the end if they cannot be relayed in a timely fashion
  - Or some may be answered by the moderators in chat directly
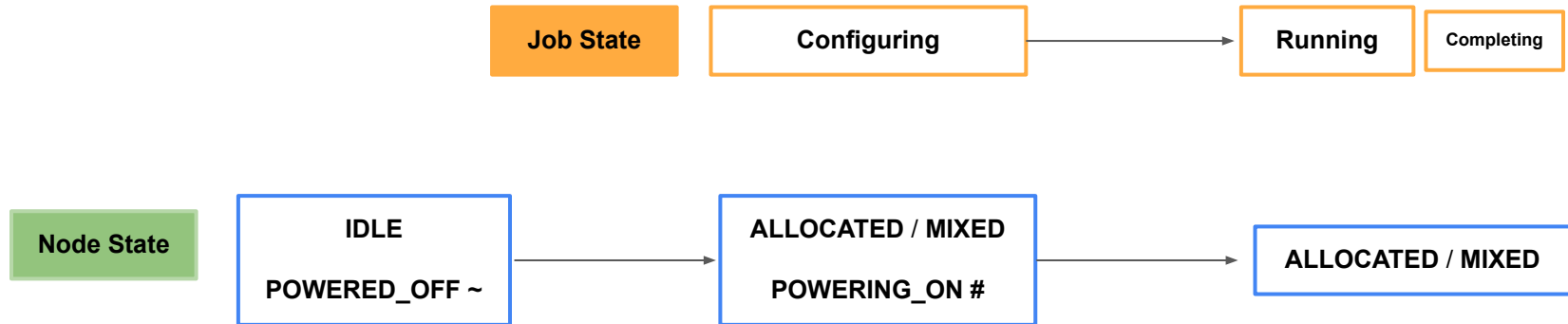
# Slurm in the Clouds

Nick Ihli
SchedMD

- New Power Save/Cloud-related features in 21.08
- Update on Slurm in public clouds
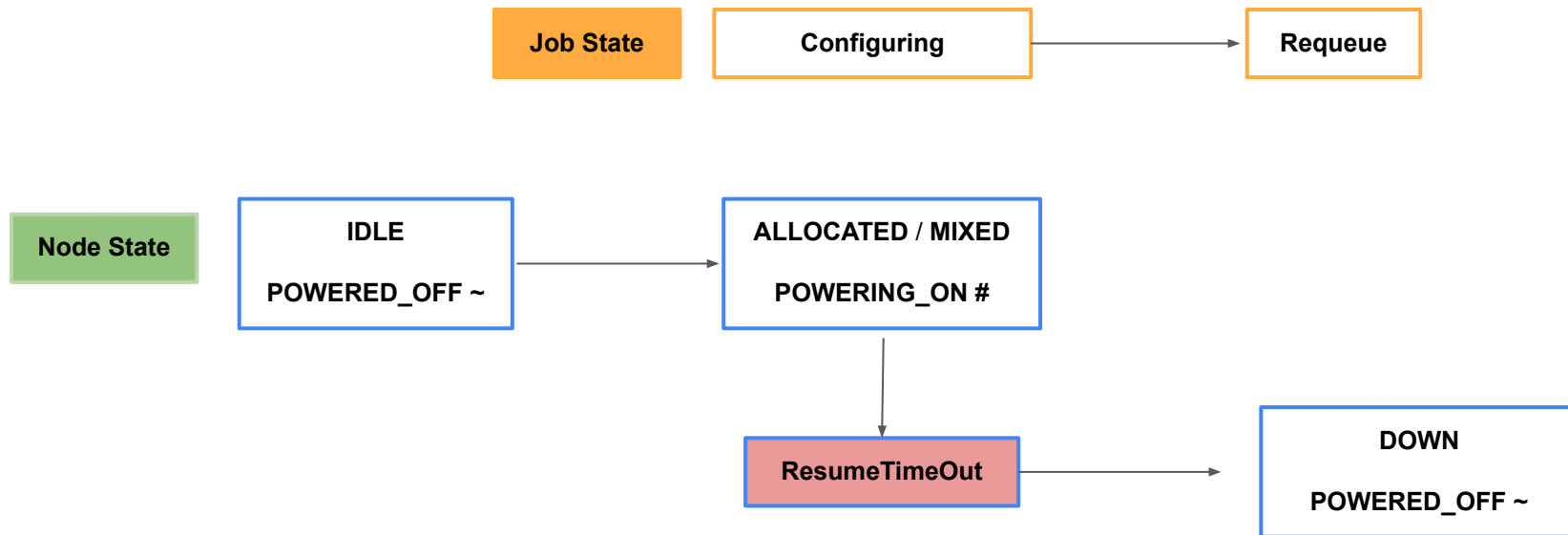
# PowerSave/Cloud changes

- SlurmctldParameters=node_reg_mem_percent
  - Allows node to register with a percentage of configured memory.
  - Defaults:
    - 90% for cloud nodes
    - 100% for everything else
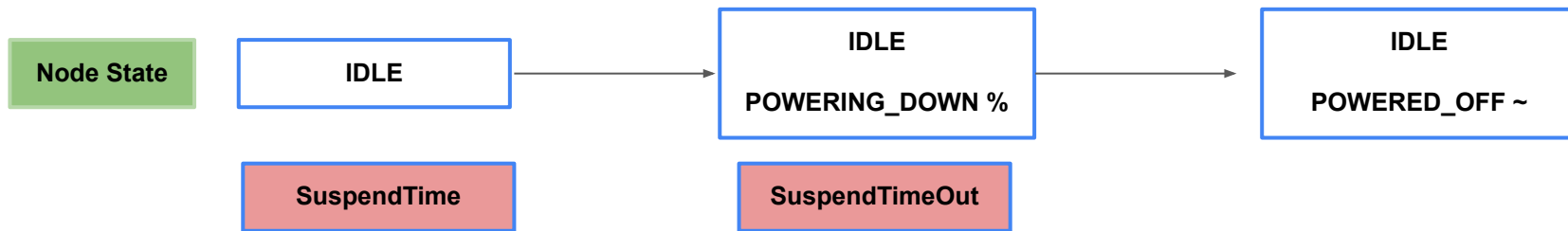
# Power State Transition - Resume

**Job State**  →  Configuring  →  Running  | Completing

**Node State**  →  IDLE

POWERED_OFF ~  →  ALLOCATED / MIXED

POWERING_ON #  →  ALLOCATED / MIXED

# Power State Transition - Resume Failure

# Power State Transition - Suspend

```
Node State        IDLE          →     IDLE                 →     IDLE

                                       POWERING_DOWN %              POWERED_OFF ~


                  SuspendTime                SuspendTimeOut
```

# PowerSave/Cloud changes

- sinfo/sview/scontrol state == base_state + flags
  - IDLE~ vs IDLE+CLOUD+POWERED_DOWN
  - sinfo -O statecomplete
  - sview "StateComplete"
  - scontrol
    - State == base_state + flags
    - Used to be shortened state + some flags
      - Old: ALLOCATED#+CLOUD
      - Now: ALLOCATED+POWERING_UP+CLOUD
  - States can be viewed through REST API

# PowerSave/Cloud changes

- scontrol update nodename=<> state=power_down[_asap|_force]
  - Like scontrol reboot states.
  - **power_down** (!) - power_down after node becomes idle
  - **power_down_asap** - put node in drain state and power down after currently running jobs
  - **power_down_force** - kill running jobs power down immediately (depends on power_save_interval)
- Able to suspend nodes that are part of SuspendExc<Parts|Nodes>

# PowerSave/Cloud changes

- LastBusy time in scontrol
  - last_busy + suspend_time < now == Suspend

# PowerSave/Cloud changes

- SuspendTime, SuspendTimeout, ResumeTimeout on partitions
  - SuspendTime can enable PowerSave if disabled at global level
  - This will be helpful for "hybrid" (i.e. bursting from on-premise) scenarios
    - e.g. by default PowerSave wants to suspend everything that isn't in SuspendExc<Nodes|Parts>
      - You have to remember to update these
    - Now with SuspendTime on the partition, you can disable PowerSave at the global level and enable on specific cloud partitions
    - e.g.
      - SuspendTime=INFINITE
      - PartitionName=cloud … SuspendTime=300

# PowerSave/Cloud changes

- JSON mapping of jobs to nodes available in ResumeProgram

```
SLURM_RESUME_FILE=/proc/1647372/fd/7:
{
    "all_nodes" : "cloud[1-3]",
    "jobs" : [
        {
            "job_id" : 140814,
            "nodes" : "cloud[1-3]",
        },
        {
            "job_id" : 140815
            "nodes" : "cloud[1-2]",
        }
    ]
}
```

# Burst Buffer

- Review Marshall's presentation from earlier on Lua Burst Buffer
- Cloud potential -  hybrid and all-in the cloud
  - Stage data before and after jobs are completed without wasting $ on compute idle time

# Cloud Partners

We have strong relationships with our public cloud partners and are working with them closely on development and consultative engagements to continue to enhance the experience of using Slurm on their clouds.

- AWS
- Microsoft Azure
- Google Cloud

# Slurm on Google Cloud

Latest Updates and Features In Development

**Open Source on SchedMD's Github:**
https://github.com/schedmd/slurm-gcp

"Version 4" Features Available Today:
- **Terraform** support generally available
- **Google Cloud HPC VM Image**-based deployment reduces deployment time to just a few minutes
- **Placement policy** support for low latency networking
- **Bulk API** reduces deployment time by reducing API calls and performing "regional capacity finding" for large deployments, up to 1,000 instances
- **Instance templates** simplify configuration and reusability
- **Cloud Marketplace** listing simplifies small Slurm cluster deployment

Google Cloud

# Schedmd-Slurm-GCP

SchedMD/Slurm

Speeding HPC, HTC & AI workloads via demand-based cloud clusters

**LAUNCH**   **VIEW PAST DEPLOYMENTS**

**OVERVIEW**   PRICING   DOCUMENTATION   SUPPORT

## Overview

Slurm is a free open-source workload manager designed specifically to satisfy the demanding needs of high performance and high throughput computing and AI. Slurm® is the industry-leading open source workload manager to manage large-scale, complex HPC and AI workloads in the cloud with faster processing and optimal consumption of the specialized resource capabilities needed for each workload. Slurm maximizes workload throughput, scale, reliability, and results in the fastest possible time, managing workloads across cloud and on-prem clusters. SchedMD® and Google have partnered to optimize Slurm for Google Cloud, enabling seamless management and integration with optimal use and scaling of Google Cloud resources so organizations can focus on getting to results faster and easier.

Learn more ↗

### About SchedMD/Slurm

SchedMD is the core developer and services provider for the market-leading Slurm workload manager. Support, consulting, configuration, development and training services accelerate workloads and results across cloud and on-

## New Schedmd-Slurm-GCP deployment

Deployment name
schedmd-slurm-gcp-1

Cluster name                                    ❓

Zone
us-central1-a                              ▼    ❓

GPU availability is limited to certain zones. Learn more ↗

⌄ **MORE**

### Network

**Network interfaces**

default default (10.128.0.0/20)                    ⌄

ADD NETWORK INTERFACE

☑ Controller External IP  ❓
Enable Private Google access or add a Cloud Router NAT on the target subnetwork before disabling

☐ Login External IP  ❓
☐ Compute Node External IPs  ❓

### Network Storage Mount

⌄ SHOW NETWORK STORAGE MOUNT OPTIONS

### Slurm Controller

**Controller Machine type**

### Schedmd-Slurm-GCP overview

Product provided by SchedMD/Slurm

### Software

| | |
|---|---|
| **Operating System** | CentOS(7) |
| **Software** | Slurm(20.11.7) |

### Documentation

Slurm Documentation ↗
Documentation for the latest release of Slurm.

### Terms of Service

By deploying the software or accessing the service you are agreeing to comply with the SchedMD/Slurm terms of service ↗, GCP Marketplace terms of service and the terms of applicable open source software licenses bundled with the software or service. Please review these terms and licenses carefully for details about any obligations you may have related to the software or service. To the limited extent an open source software license related to the software or service expressly supersedes the GCP Marketplace Terms of Service, that open source software license governs your use of that software or service.

By using this product, you understand that certain account and usage information may be shared with SchedMD/Slurm for the purposes of financial accounting, sales attribution, performance analysis, and support. ❓

Google is providing this software or service "as-is" and any support for this software or service will be provided by SchedMD/Slurm under their terms of service.

Google Cloud

# Slurm on Google Cloud

Latest Updates and Features In Development

21.08 release soon:

- **Intel Select Solution** support built-in to improve compatibility and performance

"Version 5" Features In Development:

- **Billing insights** provide visibility into Slurm and GCP billing
- **Partition flexibility** allows more flexible (re)configuration of partitions
- **Data migration** integrates Slurm data migration abilities with GCS
- **SMT configurability** integrates GCP SMT controls with Slurm-GCP
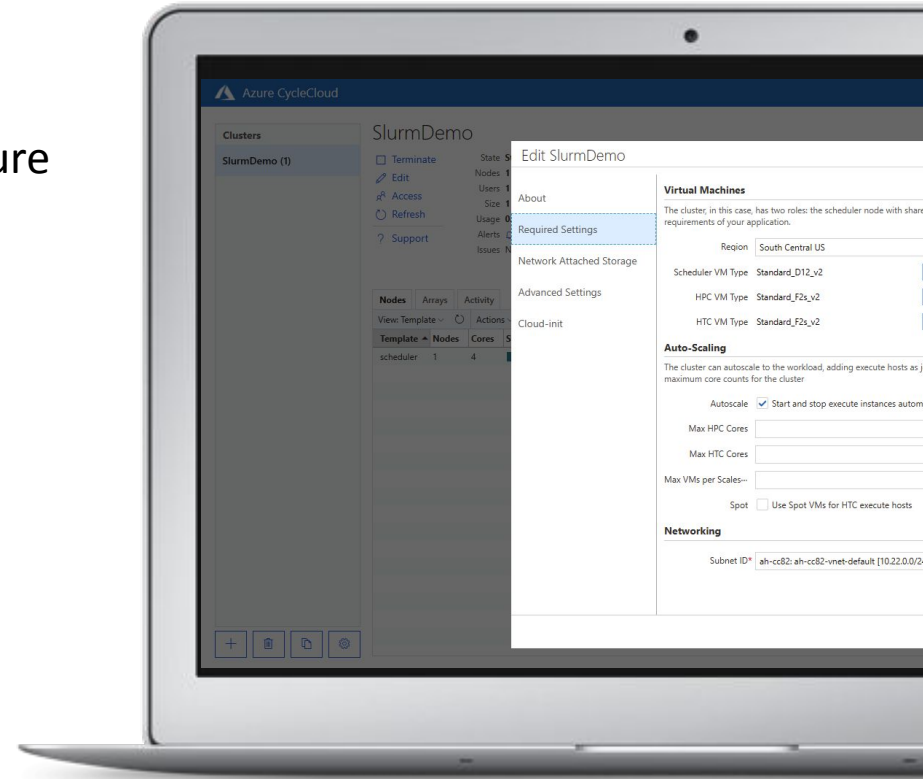- **Cloud Foundations Toolkit** simplifies and standardizes the Slurm-GCP scripts



Google Cloud

# Azure CycleCloud + Slurm

- Fastest way to get started with Slurm on Azure
  - Deploy a complete HPC cluster in just minutes
  - Azure-enable workflows with no changes
- Cluster autoscaling
  - Automatic or manual resource control
- Cost reporting and controls
  - Near real-time cost reporting
  - Link usage to spend
  - Tools to manage and control costs
- Hybrid workflows
  - Burst on-premises Slurm clusters to the cloud
- Authorization and governance:
  - AD and Azure AD integration
  - Audit and event logging
  - RBAC authorization control

Visit https://azure.microsoft.com/en-us/features/azure-cyclecloud to learn more

Microsoft Azure

# Microsoft Azure Slurm Integration

Azure CycleCloud

+ Full support for Slurm 20.11+ with CycleCloud 8.2

Support for job topology of MPI jobs requiring InfiniBand

VMs with GPUs automatically configured with GRES settings

Easy to use job accounting configuration

Burst execute nodes from on-premises to Azure

COMING SOON: Custom Slurm.conf settings set via the CycleCloud UI

**Slurm Settings**

Section for configuring Slurm

| | |
|---|---|
| Slurm Version* | 20.11.7-1 |
| Job Accounting | ☑ Configure Slurm job accounting |
| Slurm DBD URL | myslurmacctdb.mariadb.database.azure.com |
| Slurm DBD User | myacctdbuser |
| Slurm DBD Password | ●●●●●●●● |
| ShutdownPolicy | Terminate |
| Additional Slurm co··· | SuspendExcParts=htc<br>Prolog=/my/custom/prolog<br>MpiDefault=pmix |

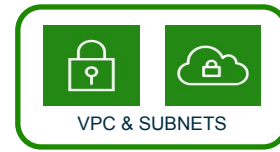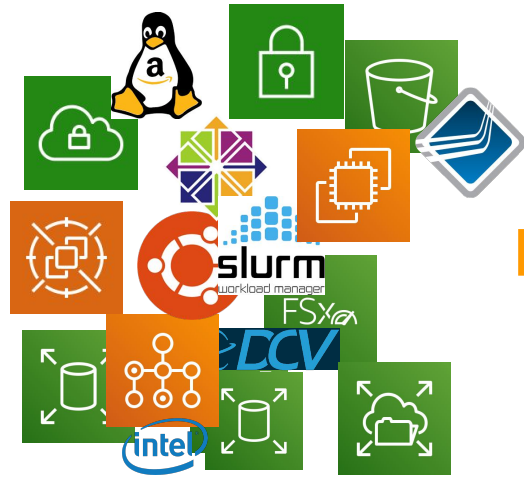Microsoft Azure    Official project repository: https://github.com/Azure/cyclecloud-slurm

# AWS and Slurm Updates

- Slurm 21.08 – Key new features
  - Improved "all or nothing" allocation and scaling
  - Support for Slurm REST API with Amazon Cognito
- ParallelCluster 3.0 – Launched Sep 10
  - AWS-supported, open source cluster management tool that is the simplest way to get started with Slurm and HPC on AWS
  - Enables:
    - Easy Cluster Management
    - Automatic Resource Scaling
    - Seamless Migration to the Cloud
  - Integration with Slurm 21.08 coming soon

# AWS ParallelCluster and Slurm



**Enable On-premises environmental parity | Facilitate "lift and shift" and applications migration over time.**

# Questions?

# Next Session

- The next presentation is by Tim Wickberg: "Slurm 21.08 and Beyond"
- Starts at 11:30am Mountain Daylight Time (UTC-6)
- And is on a separate YouTube Live stream
- Please see the SchedMD Slurm YouTube channel for links

# End Of Stream

- Thanks for watching!