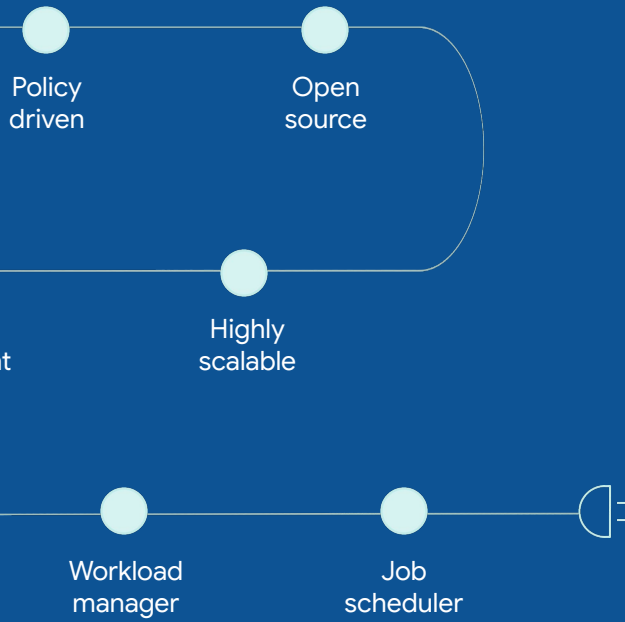


Accelerating HPC and AI with Slurm and SchedMD

Nick Ihli, Director - Solutions Engineering and Cloud
nick@schedmd.com



Most people know Slurm!



Allocates access

to resources to users for some duration of time for a workload

Provides framework

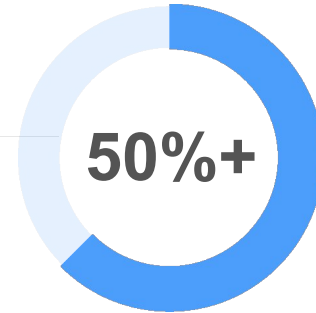
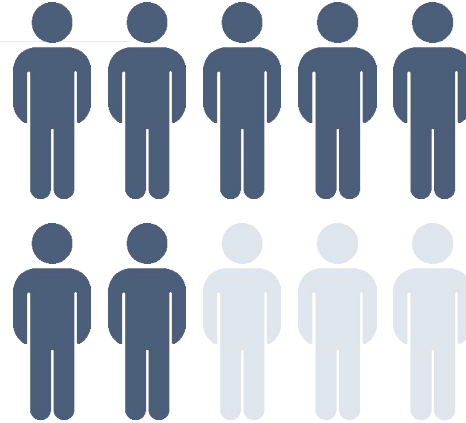
for starting, executing, and monitoring work on the set of allocated nodes

Arbitrates contention

for resources by managing a queue of pending work

Slurm on TOP500

7 of the TOP10
And more than **50% of the**
TOP500 use Slurm



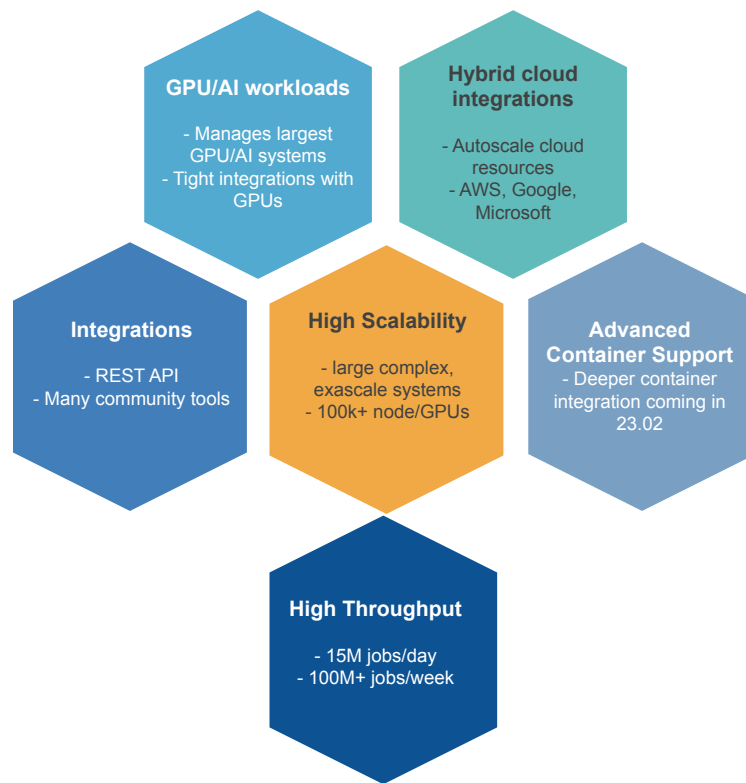
But what is SchedMD?

Maintainers and supporters of Slurm

- Only organization providing level-3 support
- Training
- Consultation
- Custom Development



Slurm leads in industry trends

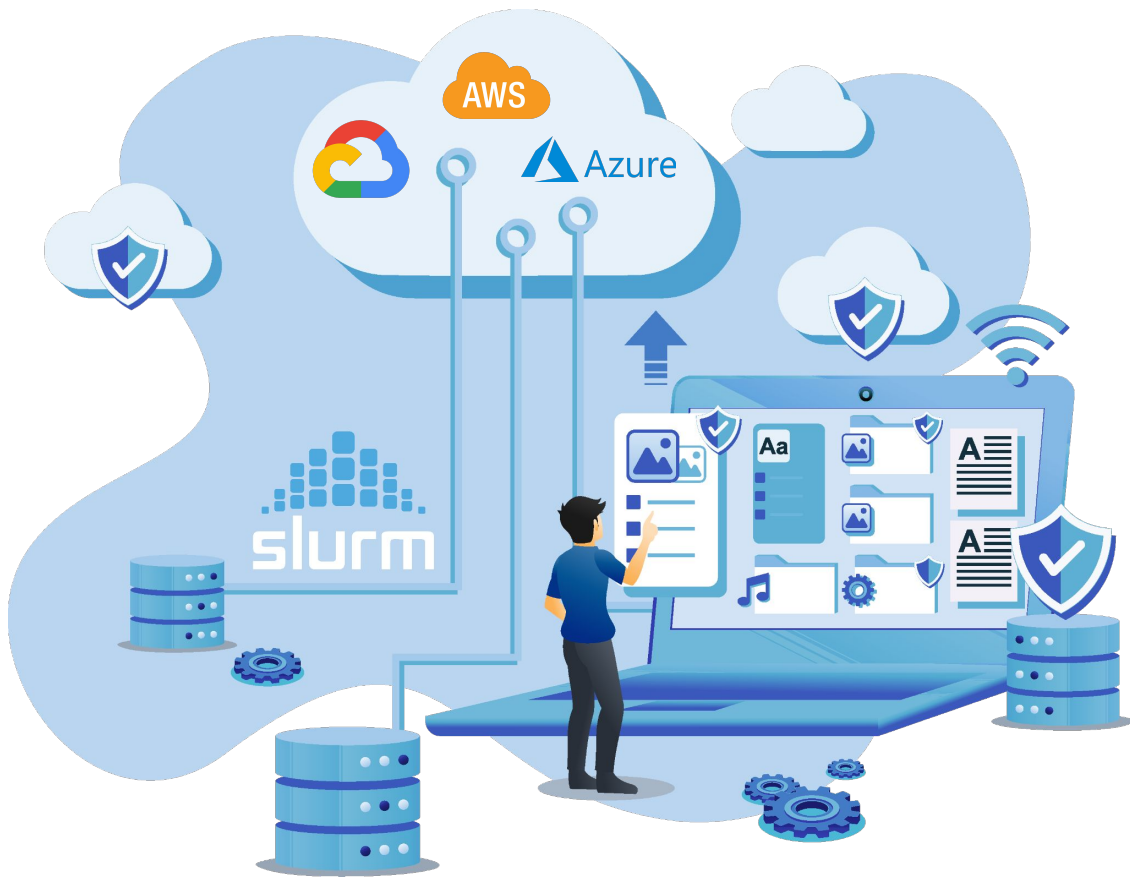


Do More with Slurm



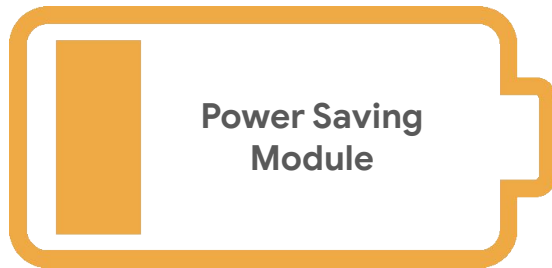
Tight GPU integration

- GPUs are a first class citizen like a CPU
- Allows for fine-grained GPU requests
- Bind tasks to GPUs
- MIG support
- Auto-detecting of GPUs
- Constrain jobs to allocated GPUs
- Sets `CUDA_VISIBLE_DEVICES` environment variable



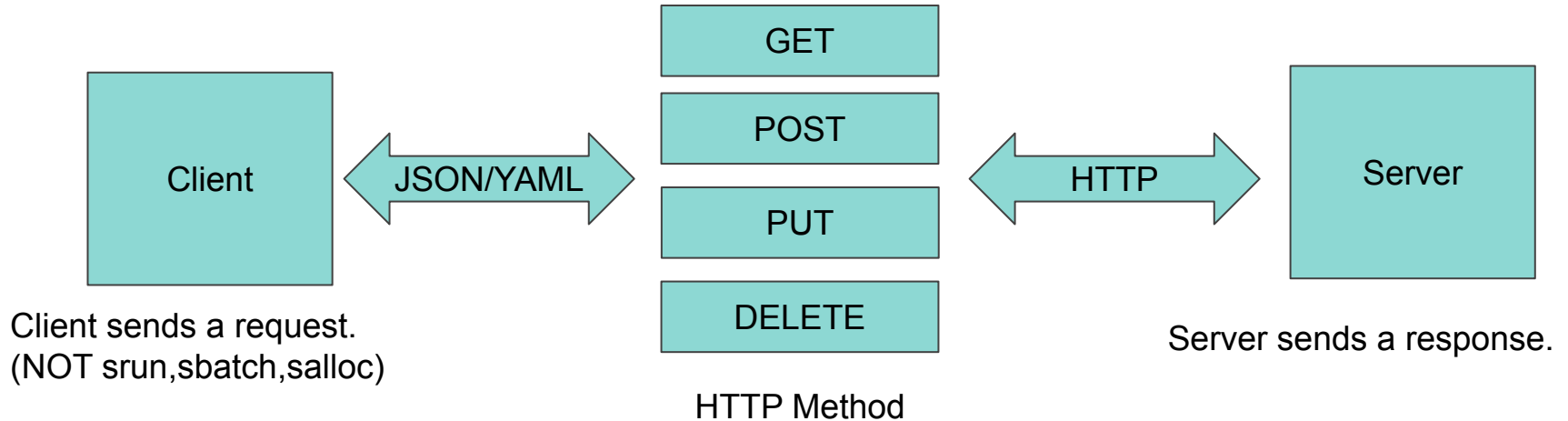
**Slurm is the bridge
between on-prem
and cloud**

Slurm Cloud Autoscaling



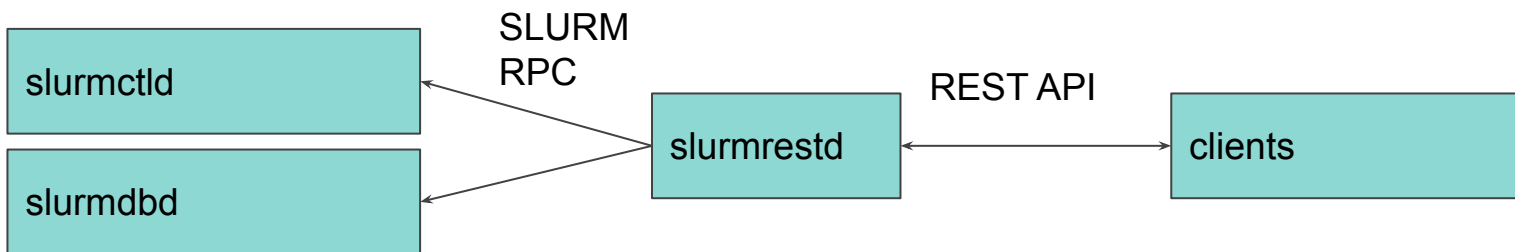
- Resume program
 - nodes are “powered on” when required
 - Suspend program
 - nodes are “powered off” when no longer required
- Suspend Time
 - Time for node to be IDLE before it is put in power saving mode (deprovisioned)
 - Set globally or per partition
 - Timeout settings
 - When to fail if the node has not registered with the controller

Application Integration - Slurm REST API



slurmrestd

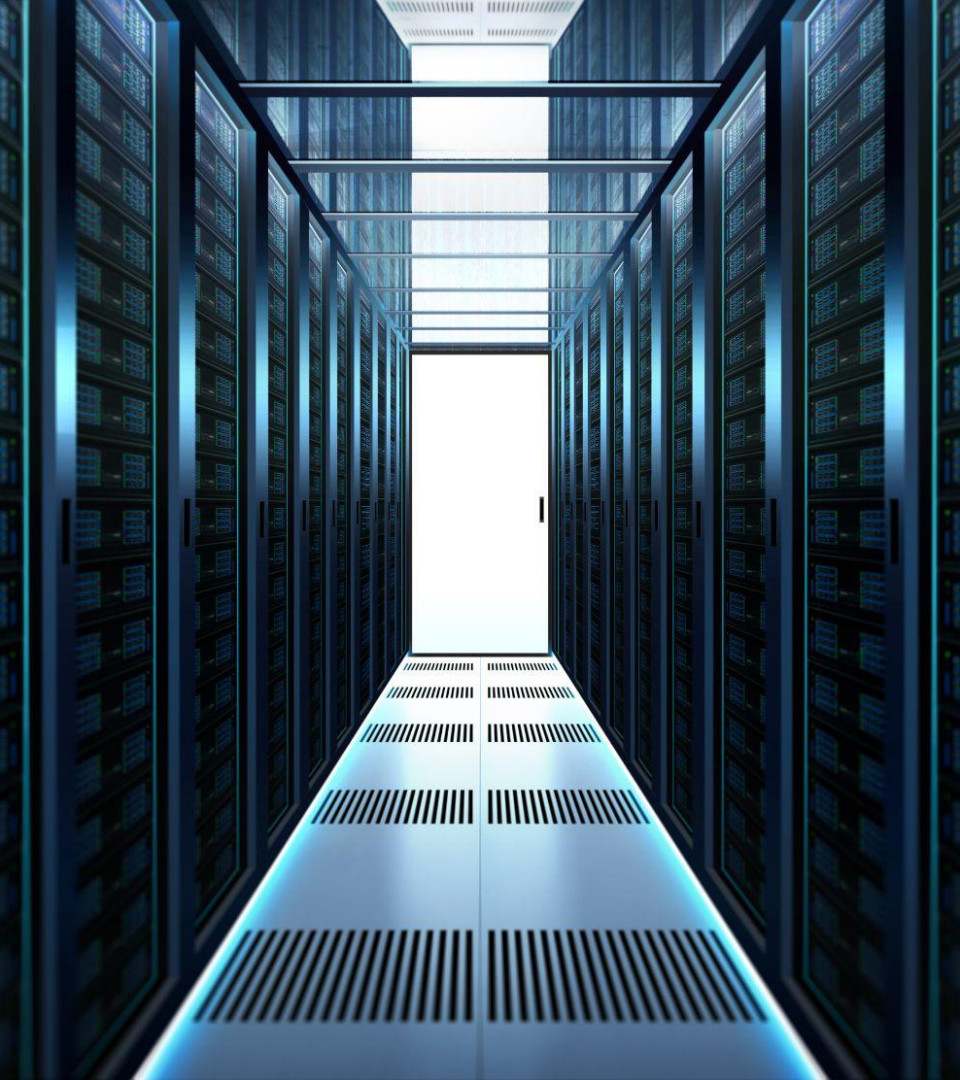
A tool that runs inside of the Slurm perimeter that will translate JSON/YAML requests into Slurm RPC requests



Tighter Container Integration - 23.02

- New scrun daemon - goal is to make containers boring for users
- Users have better things to do than learn about the intricacies of containers
- Use Slurm's existing infrastructure to run containers on compute nodes
- Automatic staging out and in of containers controlled by system administrators
- End requirement that users manually prepare their images on compute nodes.
- Interface directly with OCI runtime clients (Docker or Podman or ...)

Do More with SchedMD: A Migration Journey



A Migration Journey

Large Energy Company

- Using their scheduler for many years
 - Can't just flip a switch and go to production
- Massive scale
 - multiple international sites, nodes, and workloads
- Many integration tools required

3-4 months production

Three Migration Steps

Admin/user education

Training: Help admins identify the commonalities and learn the Slurm way

Wrappers: Use as a bridge to migration not a crutch

- LSF, Grid Engine - command and submission
- PBS - command, submission, environment variables, #PBS scripts

Policy replication

Reevaluate policies

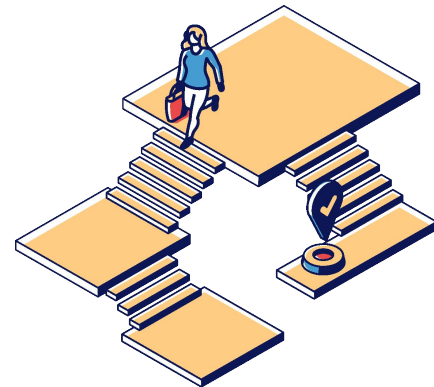
- Are we continuing to produce technical debt due to “doing things how we’ve always done them?”
- Opportunity to take a step back and redefine policies based on Slurm best practices

Optimizing for scale and throughput

Tooling integration

Most time-consuming of the migration journey.

- REST API
- Community integrations



Thank You

schedmd.com

slurm.schedmd.com

nick@schedmd.com