# Introduction to Slurm

Tim Wickberg

SchedMD

Slurm User Group Meeting 2017

# Outline

- Roles of resource manager and job scheduler
- Slurm description and design goals
- Slurm architecture and plugins
- Slurm configuration files and commands
- Accounting

# Outline

- **Roles of resource manager and job scheduler**
- Slurm description and design goals
- Slurm architecture and plugins
- Slurm config files and commands
- Accounting

# Role of a Resource Manager

- The "glue" for a parallel computer to execute parallel jobs
- It should make a parallel computer as almost easy to use as a PC
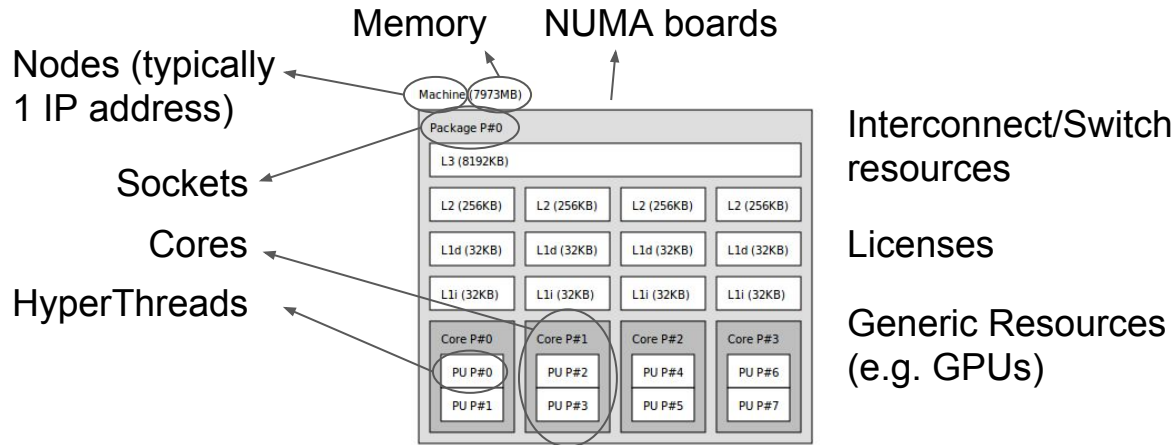
```
On a PC.
Execute program "a.out"

a.out
```

```
On a cluster.
Execute 8 copies of "a.out"

srun -n8 a.out
```

- MPI would typically be used to manage communications within the parallel program
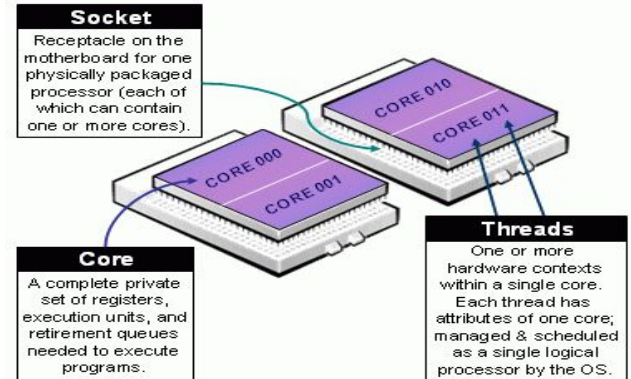
# Roles of a Resource Manager

- ## Allocate resources within a cluster

Memory    NUMA boards

Nodes (typically 1 IP address)

Machine (7973MB)

Package P#0

L3 (8192KB)

Sockets

| L2 (256KB) | L2 (256KB) | L2 (256KB) | L2 (256KB) |

| L1d (32KB) | L1d (32KB) | L1d (32KB) | L1d (32KB) |

Cores

| L1i (32KB) | L1i (32KB) | L1i (32KB) | L1i (32KB) |

HyperThreads

| Core P#0 | Core P#1 | Core P#2 | Core P#3 |
| PU P#0 | PU P#2 | PU P#4 | PU P#6 |
| PU P#1 | PU P#3 | PU P#5 | PU P#7 |

Interconnect/Switch resources

Licenses

Generic Resources (e.g. GPUs)

**Socket**
Receptacle on the motherboard for one physically packaged processor (each of which can contain one or more cores).

CORE 010
CORE 011
CORE 000
CORE 001

**Core**
A complete private set of registers, execution units, and retirement queues needed to execute programs.

**Threads**
One or more hardware contexts within a single core. Each thread has attributes of one core; managed & scheduled as a single logical processor by the OS.

- ## Launch and otherwise manage jobs

# Role of a Job Scheduler

- When there is more work than resources, the job scheduler manages queue(s) of work
    - Supports complex scheduling algorithms
        - Optimized for network topology, fair-share scheduling, advanced reservations, preemption, gang scheduling (time-slicing jobs), backfill scheduling, etc.
        - Job can be prioritized using highly configurable parameters such as job age, job partition, job size, job QOS, etc.
    - Supports resource limits (by queue, user, group, etc.)

# Examples

| Resource Managers | Schedulers |
|---|---|
| ALPS (Cray) | Maui |
| Torque | Moab |
| LoadLeveler (IBM) | |
| Slurm | |
| LSF | |
| PBS Pro | |

Many span both roles

Slurm started as a resource manager (the "rm" in Slurm) and added scheduling logic later

# Outline

- Roles of resource manager and job scheduler
- **Slurm description and design goals**
- Slurm architecture and plugins
- Slurm config files and commands
- Accounting

# What is Slurm?

- Historically Slurm was an acronym standing for
  - **S**imple **L**inux **U**tility for **R**esource **M**anagement
- Development started in 2002 at Lawrence Livermore National Laboratory as a resource manager for Linux clusters
- Sophisticated scheduling plugins added in 2008
- About 500,000 lines of C code today (plus test suite and docs)
- Used on many of the world's largest computers
- Active global development community

# Slurm Design Goals

- Highly scalable (managing 3.1 million core Tianhe-2, tested to much larger systems using emulation)
- Open source (GPL version 2, available on Github)
- System administrator friendly
- Secure
- Fault-tolerant (no single point of failure)
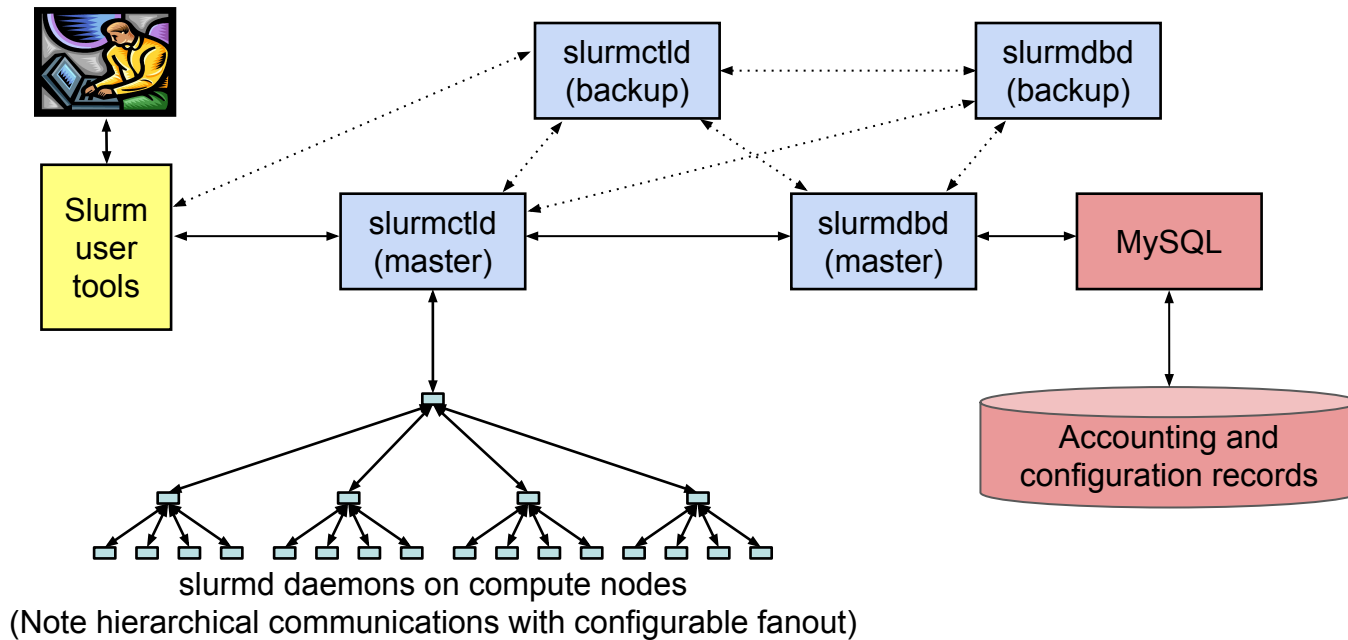- Portable - targeting POSIX2008.1 and C99

# Slurm Portability

- *Autoconf* configuration engine adapts to environment
- Provides scheduling framework with general-purpose plugin mechanism. System administrator can extensively customize installation using a building- block approach
- Various system-specific plugins available and more under development (e.g. *select/bluegene*, *select/cray*)
- Huge range of use cases:
  - Intel's "cluster on a chip": Simple resource manager
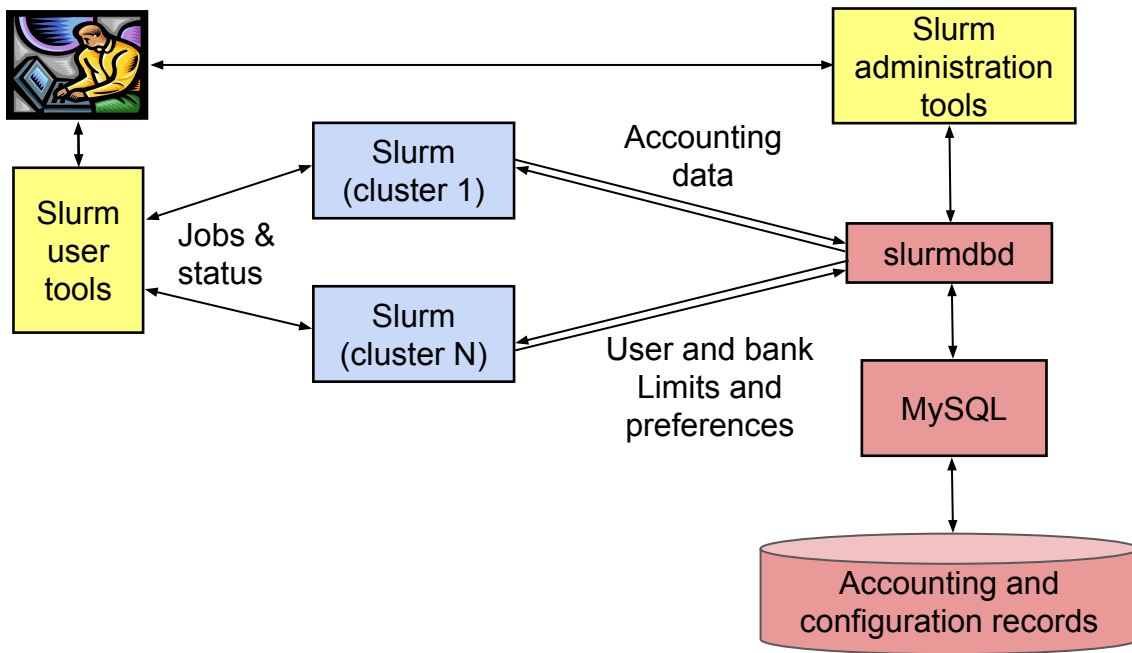  - Sophisticated workload management at HPC sites

# Outline

- Roles of resource manager and job scheduler
- Slurm description and design goals
- **Slurm architecture and plugins**
- Slurm config files and commands
- Accounting

# Cluster Architecture



slurmd daemons on compute nodes
(Note hierarchical communications with configurable fanout)

Copyright 2017 SchedMD LLC
http://www.schedmd.com

# Typical Enterprise Architecture



Copyright 2017 SchedMD LLC
http://www.schedmd.com

# Daemons

- **slurmctld** – Central controller (typically one per cluster)
  - Monitors state of resources
  - Manages job queues
  - Allocates resources
- **slurmdbd** – Database daemon (typically one per enterprise)
  - Collects accounting information
  - Uploads configuration information (limits, fair-share, etc.) to slurmctld

# Daemons

- **slurmd** – Compute node daemon (typically one per compute node)
    - Launches and manages slurmstepd (see below)
    - Small and very light-weight
    - Quiescent after launch except for optional accounting
    - Supports hierarchical communications with configurable fanout
- **slurmstepd** – Job step shepherd
    - Launched for batch job and each job step
    - Launches user application tasks
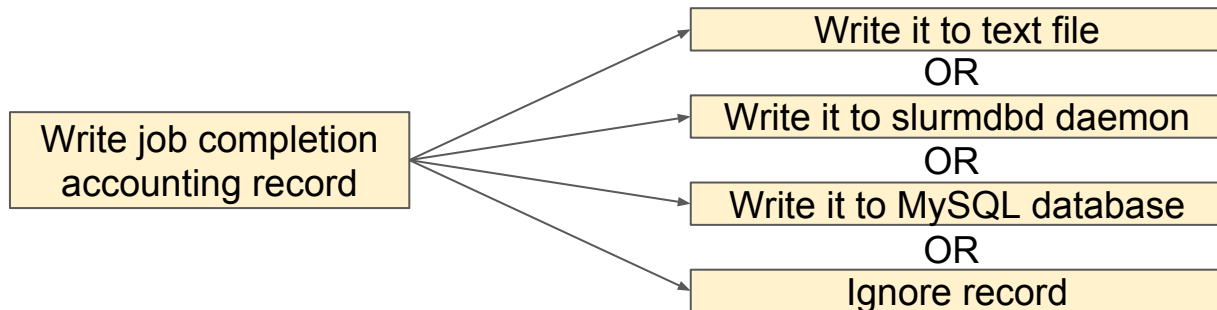    - Manages application I/O, signals, etc.

# Plugins

- Dynamically linked objects loaded at run time based upon configuration file and/or user options
- 100+ plugins of 26 different varieties currently available
  - Network topology: 3D torus, tree, etc
  - MPI: OpenMPI, MPICH1, MVAPICH, MPICH2, etc
  - External sensors: Temperature, power consumption, etc.

| Slurm Kernel (65% of code) | | | | |
|---|---|---|---|---|
| Authentication Plugin | MPI Plugin | Checkpoint Plugin | Topology Plugin | Accounting Storage Plugin |
| Munge | pmi2 | BLCR | Tree | MySQL |

# Plugin Design

- Plugins typically loaded when the daemon or command starts and persist indefinitely
- Provide a level of indirection to a configurable underlying function

```
                                        ┌─────────────────────────────┐
                                        │     Write it to text file   │
                                        └─────────────────────────────┘
                                                      OR
┌─────────────────────┐                 ┌─────────────────────────────┐
│ Write job completion │───────────────▶│  Write it to slurmdbd daemon│
│  accounting record   │                └─────────────────────────────┘
└─────────────────────┘                               OR
                                        ┌─────────────────────────────┐
                                        │  Write it to MySQL database │
                                        └─────────────────────────────┘
                                                      OR
                                        ┌─────────────────────────────┐
                                        │        Ignore record        │
                                        └─────────────────────────────┘
```

# Plugin Development

- APIs are all documented for custom development (e.g. GreenSpot for optimized use of green energy sources)
- Most plugins have several examples available
- Some plugins have a LUA script interface

# Job Submit Plugin

- Call for each job submission or modification
- Can be used to set default values or enforce limits using functionality outside of Slurm proper

Two functions need to be supplied:

int job_submit(struct job_descriptor *job_desc, uint32_t submit_uid);
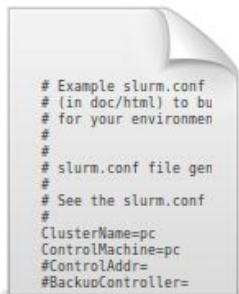int job_modify(struct job_descriptor *job_desc, struct job_record *job_ptr);

# Outline

- Roles of resource manager and job scheduler
- Slurm description and design goals
- Slurm architecture and plugins
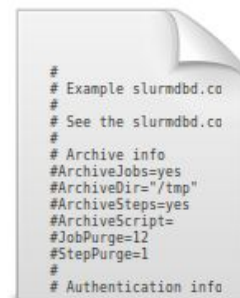- **Slurm config files and commands**
- Accounting

# Slurm Configuration

**slurm.conf**

```
# Example slurm.conf
# (in doc/html) to bu
# for your environmen
#
# slurm.conf file gen
#
# See the slurm.conf
#
ClusterName=pc
ControlMachine=pc
#ControlAddr=
#BackupController=
```

- General conf
- Plugin activation
- Sched params
- Node definition
- Partition conf

**slurmdbd.conf**

```
#
# Example slurmdbd.co
#
# See the slurmdbd.co
#
# Archive info
#ArchiveJobs=yes
#ArchiveDir="/tmp"
#ArchiveSteps=yes
#ArchiveScript=
#JobPurge=12
#StepPurge=1
#
# Authentication info
```

- Describes slurmdbd
- Archive/Purge parameters
- Storage options

# Slurm Configuration

**topology.conf**

```
# topology.conf
# Switch Configuratio
#
# Haswell
SwitchName=hsw1 Nodes
SwitchName=hsw2 Nodes
#
# Sandybridge
SwitchName=snb1 Nodes
SwitchName=snb2 Nodes
#SwitchName=snb3 Node
SwitchName=snb3 Nodes
#
# Westmere
```
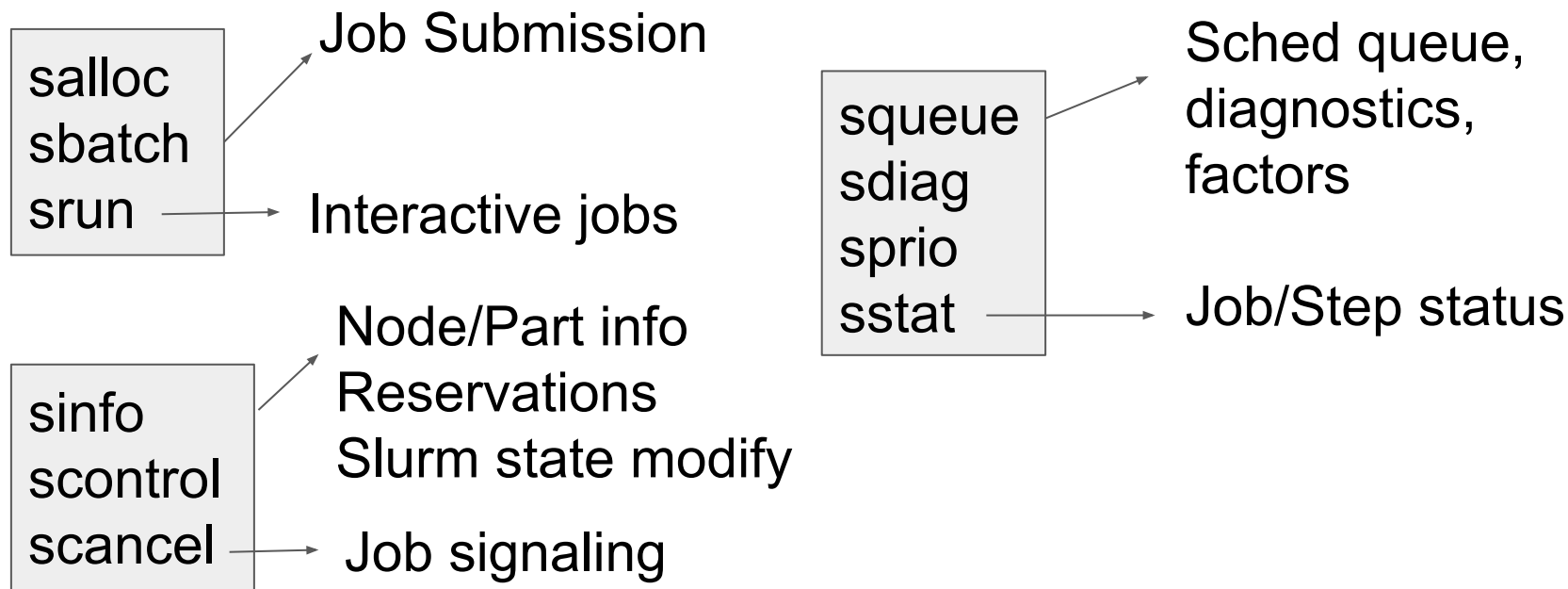
**gres.conf**

```
NodeName=compute1 Nam
NodeName=compute1 Nam
NodeName=compute2 Nam
NodeName=compute2 Nam
#NodeName=compute[1-2
```

**cgroup.conf**

```
###
#
# Slurm cgroup suppor
#
# See man slurm.conf
# information on cgro
#--
CgroupMountpoint="/sy
CgroupAutomount=yes
CgroupReleaseAgentDir
#AllowedDevicesFile="
ConstrainCores=yes
TaskAffinity=yes
ConstrainRAMSpace=yes
```

- Others: burst_buffer.conf, acct_gather.conf, knl.conf, etc.

# Commands Overview

salloc
sbatch
srun

→ Job Submission

→ Interactive jobs

sinfo
scontrol
scancel

→ Node/Part info
Reservations
Slurm state modify

→ Job signaling

squeue
sdiag
sprio
sstat

→ Sched queue,
diagnostics,
factors

→ Job/Step status

# Commands Overview

| | |
|---|---|
| sacct<br>sacctmgr<br>sshare<br>sreport | → Accounting data<br>view/modify<br>FairShare info<br>Report generation |

| | |
|---|---|
| sattach<br>sbcast<br>strigger | ↗ I/O attach to jobs,<br>file transmission<br>to nodes, events<br>triggering |

- --help, --usage
- man pages
- APIs make new tools development easier

| | |
|---|---|
| sview<br>smap | ↗ Graphical interfaces |

Copyright 2017 SchedMD LLC
http://www.schedmd.com

# Outline

- Roles of resource manager and job scheduler
- Slurm description and design goals
- Slurm architecture and plugins
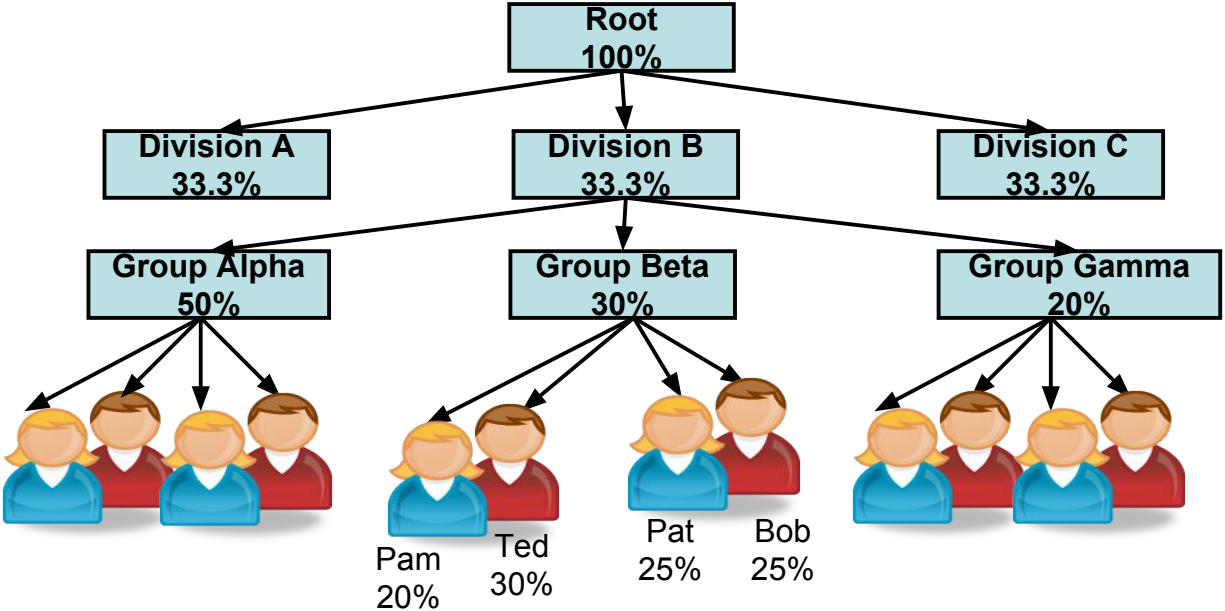- Slurm config files and commands
- Accounting

# Database Use

- Accounting information written to a database plus
  - Information pushed out live to scheduler daemons
  - Quality of Service (QOS) definitions
  - Fair-share resource allocations
  - Many limits (max job count, max job size, etc)
  - Based upon hierarchical accounts
    - Limits by user AND by accounts

*"All I can say is wow – this is the most flexible, useful scheduling tool I've ever run across."*
Adam Todorski, Rensselaer Polytechnic Institute

# Hierarchichal Account Example

# Hierarchical Accounts

- All users are not created equal
  - Different shares of resources
  - Different measures of being over- or under-served
  - Different limits
- There are many limits available
  - Per Job limits (e.g. MaxNodes)
  - Aggregate limits by user, account or QOS (e.g. GrpJobs)
  - A single user may have different shares and limits in different accounts, QOS or partitions

# Summary

- Brief introduction to Slurm
- Many more features
  - Job dependencies
  - Fine-grained task layout
  - Wrapper scripts for other workload manager command line interfaces
  - Burst Buffers, TRES, KNL support, cloud bursting, X11 forwarding, etc.
- Documentation - https://slurm.schedmd.com
- Github - https://github.com/SchedMD/slurm