# SLURM Operation on Cray XT and XE

Morris Jette
jette@schedmd.com

SchedMD LLC

# Contributors and Collaborators

This work was supported by the Oak Ridge National Laboratory Extreme Scale Systems Center.

Swiss National Supercomputing Centre performed some of the development and testing.

Cray helped with integration and testing.

# Outline

- Cray hardware and software architecture

- SLURM architecture for Cray

- SLURM configuration and use

- Status

# Cray Architecture

- Many of the most powerful computers built by Cray

- Nodes are diskless

- 2 or 3-dimension torus interconnect

  - Multiple nodes at each coordinate on some systems

- Full Linux on front-end nodes

- Lightweight Linux kernel on compute nodes
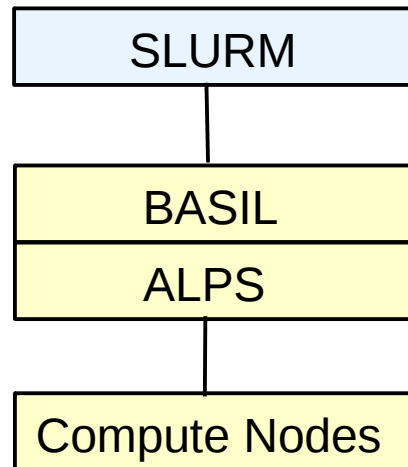
- Whole nodes must be allocated to jobs

# ALPS and BASIL

- ALPS – **A**pplication **L**evel **P**lacement **S**cheduler

    - Cray's resource manager

    - Six daemons plus variety of tools

        – One daemon runs on each compute node to launch user tasks

        – Other daemons run on service nodes

    - Rudimentary scheduling software

        – Dependent upon external scheduler (e.g. SLURM, etc) for workload management

- BASIL – **B**atch **A**pplication **S**cheduler **I**nterface **L**ayer

    - XML interface to ALPS

# SLURM Architecture for Cray

- Many tools dependent upon ALPS

    - Use SLURM as scheduler layer <u>above</u> ALPS and BASIL, not a replacement

```
┌─────────────────────┐
│        SLURM        │
└──────────┬──────────┘
           │
┌──────────┴──────────┐
│        BASIL        │
├─────────────────────┤
│        ALPS         │
└──────────┬──────────┘
           │
┌──────────┴──────────┐
│    Compute Nodes    │
└─────────────────────┘
```
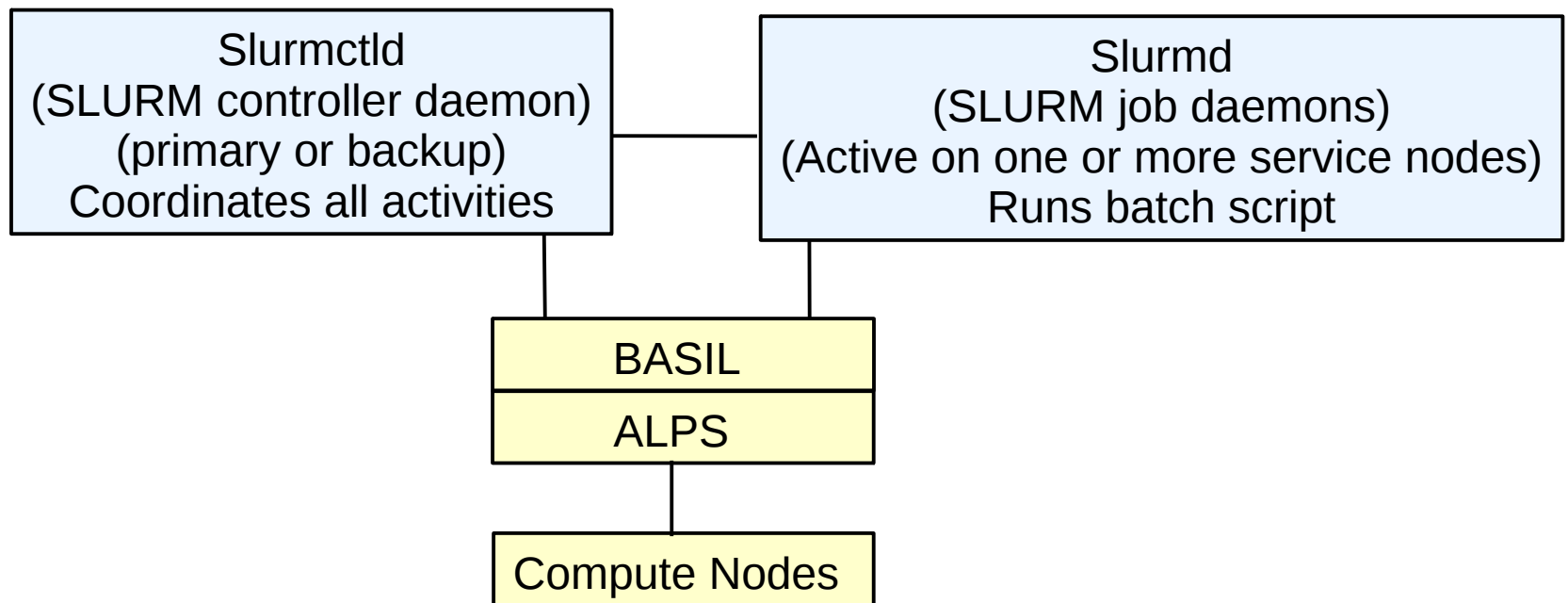
# SLURM and ALPS Functionality

- SLURM

  - Prioritizes queue(s) of work and enforces limits

  - Decides when and where to start jobs

  - Terminates job when appropriate

  - Accounts for jobs

- ALPS

  - Allocates and releases resources for jobs

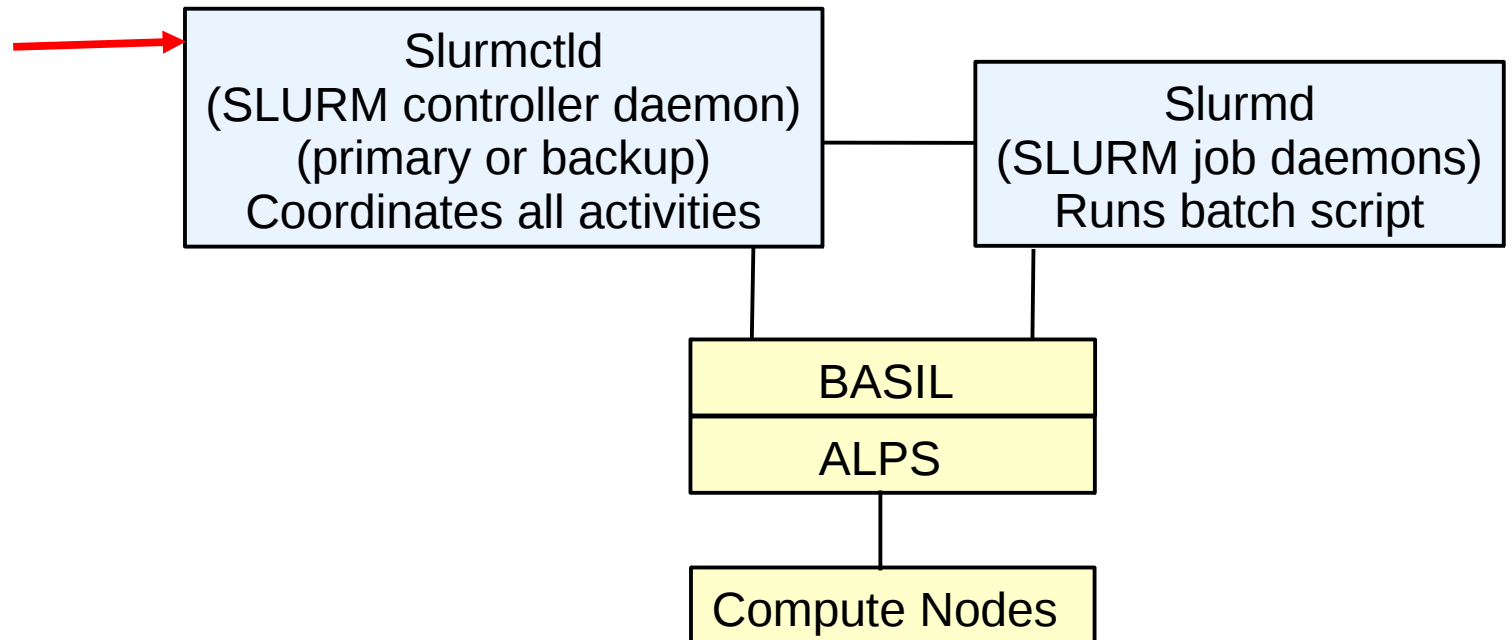  - Launches tasks

  - Monitors node health
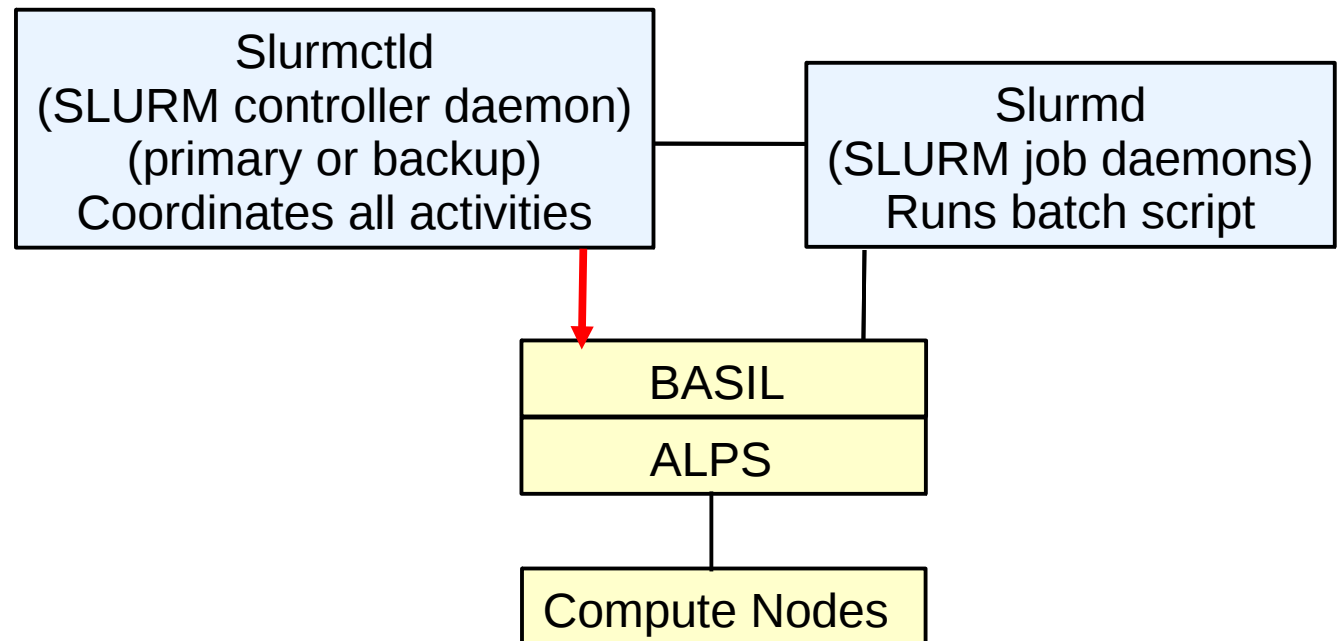
# SLURM Architecture for Cray
## (Detailed)

```
┌─────────────────────────────────┐          ┌──────────────────────────────────────────┐
│            Slurmctld            │          │                  Slurmd                  │
│   (SLURM controller daemon)     │──────────│          (SLURM job daemons)             │
│       (primary or backup)       │          │   (Active on one or more service nodes)  │
│    Coordinates all activities   │          │            Runs batch script             │
└─────────────────────────────────┘          └──────────────────────────────────────────┘
                    │                                          │
                    │         ┌──────────────────────┐         │
                    └─────────┤        BASIL         ├─────────┘
                              ├──────────────────────┤
                              │        ALPS          │
                              └──────────────────────┘
                                         │
                              ┌──────────────────────┐
                              │    Compute Nodes     │
                              └──────────────────────┘
```

# Job Launch Process

1. User submits script →

| Slurmctld |
|---|
| (SLURM controller daemon) |
| (primary or backup) |
| Coordinates all activities |

| Slurmd |
|---|
| (SLURM job daemons) |
| Runs batch script |

| BASIL |
|---|
| ALPS |

| Compute Nodes |
|---|

# Job Launch Process

1. User submits script

2. Slurmctld creates
ALPS reservation

```
┌─────────────────────────────┐        ┌──────────────────────┐
│         Slurmctld           │        │       Slurmd         │
│  (SLURM controller daemon)  │────────│ (SLURM job daemons)  │
│    (primary or backup)      │        │  Runs batch script   │
│  Coordinates all activities │        │                      │
└─────────────────────────────┘        └──────────────────────┘
              │                                    │
              ▼                                    │
        ┌──────────────────────┐                  │
        │       BASIL          │──────────────────┘
        ├──────────────────────┤
        │       ALPS           │
        └──────────────────────┘
                   │
        ┌──────────────────────┐
        │    Compute Nodes     │
        └──────────────────────┘
```

# Job Launch Process

1. User submits script

2. Slurmctld creates ALPS reservation

3. Slurmctld sends script to slurmd

```
┌─────────────────────────────┐          ┌─────────────────────────┐
│         Slurmctld           │─────────▶│        Slurmd           │
│ (SLURM controller daemon)   │          │  (SLURM job daemons)    │
│   (primary or backup)       │          │    Runs batch script    │
│  Coordinates all activities │          │                         │
└─────────────────────────────┘          └─────────────────────────┘
                    │                                  │
              ┌─────────────────────────┐
              │         BASIL           │
              ├─────────────────────────┤
              │         ALPS            │
              └─────────────────────────┘
                           │
              ┌─────────────────────────┐
              │     Compute Nodes       │
              └─────────────────────────┘
```

# Job Launch Process

1. User submits script

2. Slurmctld creates ALPS reservation

3. Slurmctld sends script to slurmd

4. Slurmd claims reservation for specific session ID and launches interpreter for script

Slurmctld
(SLURM controller daemon)
(primary or backup)
Coordinates all activities

Slurmd
(SLURM job daemons)
Runs batch script

BASIL

ALPS

#!/bin/bash
srun a.out

Compute Nodes

# Job Launch Process

1. User submits script

2. Slurmctld creates ALPS reservation

3. Slurmctld sends script to slurmd

4. Slurmd claims reservation for specific session ID and launches interpreter for script

5. aprun (optionally using the srun wrapper) launches tasks on compute nodes

| Slurmctld (SLURM controller daemon) (primary or backup) Coordinates all activities |
|---|

| Slurmd (SLURM job daemons) Runs batch script |
|---|

| BASIL |
|---|
| ALPS |

| #!/bin/bash srun a.out |
|---|

| Compute Nodes |
|---|

| aprun a.out |
|---|

# SLURM Architecture for Cray

- Almost all Cray-specific logic is in a Resource Selection plugin (as SLURM does for IBM BlueGene systems)

- The *select/cray* plugin in-turn calls the *select/linear* plugin to provide full-node resource allocation support including job preemption, memory allocation, optimized topology layout, etc.

| SLURM Kernel | | | | |
|---|---|---|---|---|
| Resource Selection Plugin | | | | Other Plugins |
| BlueGene | ConsRes | Cray | Linear | |
| | | Linear | | |

New for Cray

# SLURM's *select/cray* plugin



SchedMD LLC
http://www.schedmd.com

# SLURM Configuration

```
#
# Sample slurm.conf file for Cray system
# Selected portions
#
SelectType=select/cray          # Communicates with ALPS
#
FrontEndName=front[00-03]       # Where slurmd daemons run
NodeName=nid[00000-01023]
PartitionName=debug Nodes=nid[00000-00015] MaxTime=30
PartitionName=batch Nodes=nid[00016-01023] MaxTime=24:00:00
```

# *srun* Command

- SLURM has *srun* command for Cray systems that is a wrapper for both *salloc* (to allocate resources as needed) and *aprun* (to launch tasks)

- Options are translated to the extent possible

- Build SLURM with *configure* option *–with-srun2aprun* to build wrapper

  - Otherwise *srun* command advises use of *aprun* and exits

# *srun* Options

- If no allocation exists, the *srun* options are translated directly to *salloc* options to create a job allocation

  - Many srun option only apply when a job allocation is created

- After an allocation is created

  - Most common options are translated from *srun* to *aprun* (task count, node count, time limit, file names, support for multiple executables, etc.)

  - Some *srun* options lack *aprun* equivalent and vice-versa

  - srun's "*–alps=*" option can pass any other options to *aprun*

- SLURM environment variables are not currently set

- There are fundamental differences in I/O

  - For example, ALPS does not support per-rank I/O streams

# *sview* of Emulated System

# *sview* of Emulated System



SchedMD LLC
http://www.schedmd.com

# *smap* of Emulated System

# Caveats

- Some SLURM functionality has no ALPS equivalent

  - Independent I/O by task

  - Output labeled by task ID

- Some ALPS options have no SLURM equivalent

  - *srun* wrapper has *–alps* option to pass arbitrary arguments to *aprun*

- Some options are similar, but impossible to directly translate

  - Task binding syntax

  - Per-task vs. per-CPU limits

# Caveats
## (continued)

- SLURM environment variables are not set
  - Many need to be set on a per-node or per-task basis, so ALPS must do this
  - Under development by Cray

# Caveats
## (continued)

- SLURM GUIs (actually the *curses* and GTK libraries they use) have limited scalability

  - Scales to a few thousand nodes

  - Currently each position displayed represents a unique X, Y, Z-coordinate

  - If multiple nodes share an X, Y, Z-coordinate, the information for only one node is displayed

  - We found this better than displaying each node independently and providing confusing topology information, but could change this if desired

# Status

- SLURM has been running reliably over ALPS at Swiss National Supercomputer Centre since April 2011

- Validated by Cray in July/August 2011