





Software accounting using SGAS, Slurm & CGroups

Magnus Jonsson

<magnus@hpc2n.umu.se>

<http://www.hpc2n.umu.se/>





Background

- What software are our users running?
- What version of the software are they using?
- How much of the software are they using?
- Are the users running the software they say they are using?
- Are we putting optimization effort on the right software?





Background

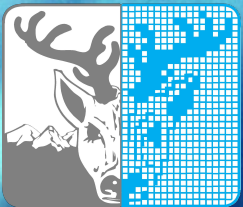
- This must be doable?





Background

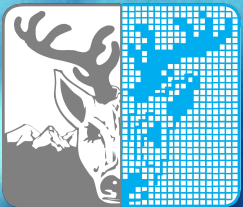
- This must be doable?
- Or?





Background

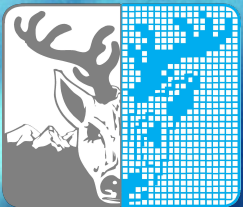
- This must be doable?
- Or?
- Proof of concept





Techniques used/prerequisite

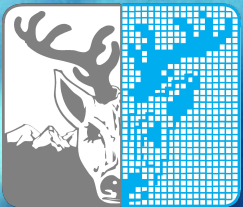
- /proc/
- Slurm
- CGroups
- perl
- json
- SGAS





CGroups

- Linux
- Contain
 - CPU
 - Memory





Slurm

- Uses CGroups to contain jobs
 - Memory
 - CPU
 - ...
- One CGroup for each job per node





/proc/

- CGroups leaves a trace in /proc
- /proc/<pid>/cpuset
 - /slurm/uid_1234/job_1234567/step_1
- /proc/<pid>/exe
 - Symbolic link to the binary running
- /proc/<pid>/task/stat
 - Same as /proc/<pid>/stat but for every thread.
 - User/System/"walltime"/vsize/RSS





Collector

- Perl script that runs on every node
 - Daemon
 - Cron
- Takes $< 0.1s$ to run
 - Including loading perl and modules





Problems?

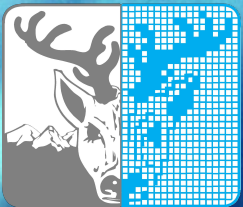
- User compiled software
 - Do we care?
- When is a job completed?
 - Assume
 - Ask Slurm
- What to save?
- Sample
 - What are we missing?
- Translate running exec to software





When is a job completed?

- Assume?
- Ask slurm?





What to save?

- To much to save it all... or?
- Aggregated user/system time for each executable
- Some other metadata
- Save the largest 90% by cputime





Translate executable to program

- Orders list
- Regular expression to map path to "software" + "version"
 - Software/Version can be fixed or collected as part of path
 - User provided



How to save?

- Currently saves files into our // file system
 - Lots of files...
- Plugin for SGAS to save info in the same database.
 - Information about jobs already available in the same database
 - Easy to do with the new plugin infrastructure





Sample problem

- Today
 - Cron every 5min
- Next generation
 - Cron every 5min
 - Daemon every minute
- Short running jobs/processes
- Kernel support?
- LD_PRELOAD?





Post Processing

- Aggregate software for each job
 - Read Slurm accounting database
 - Project
 - Number of nodes
 - Job Walltime
 - ...?
 - Read from SGAS





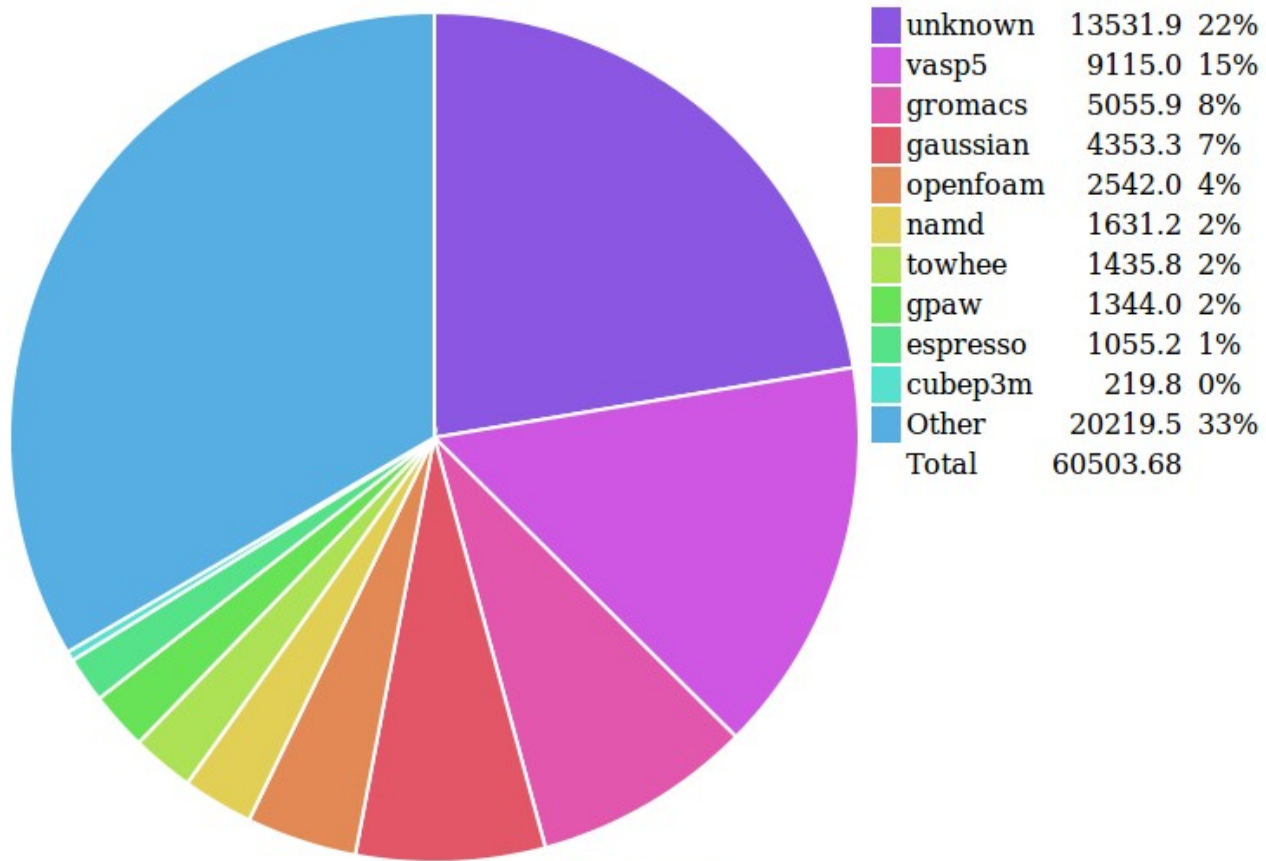
Make pretty graphs

- We are using a JavaScript tool.
 - Developed by me for an private project
 - All calculations on client side
 - Can be made into a single html file
 - Usable limit around 200-300k jobs



Show pretty graphs...

Time - Program (top10)



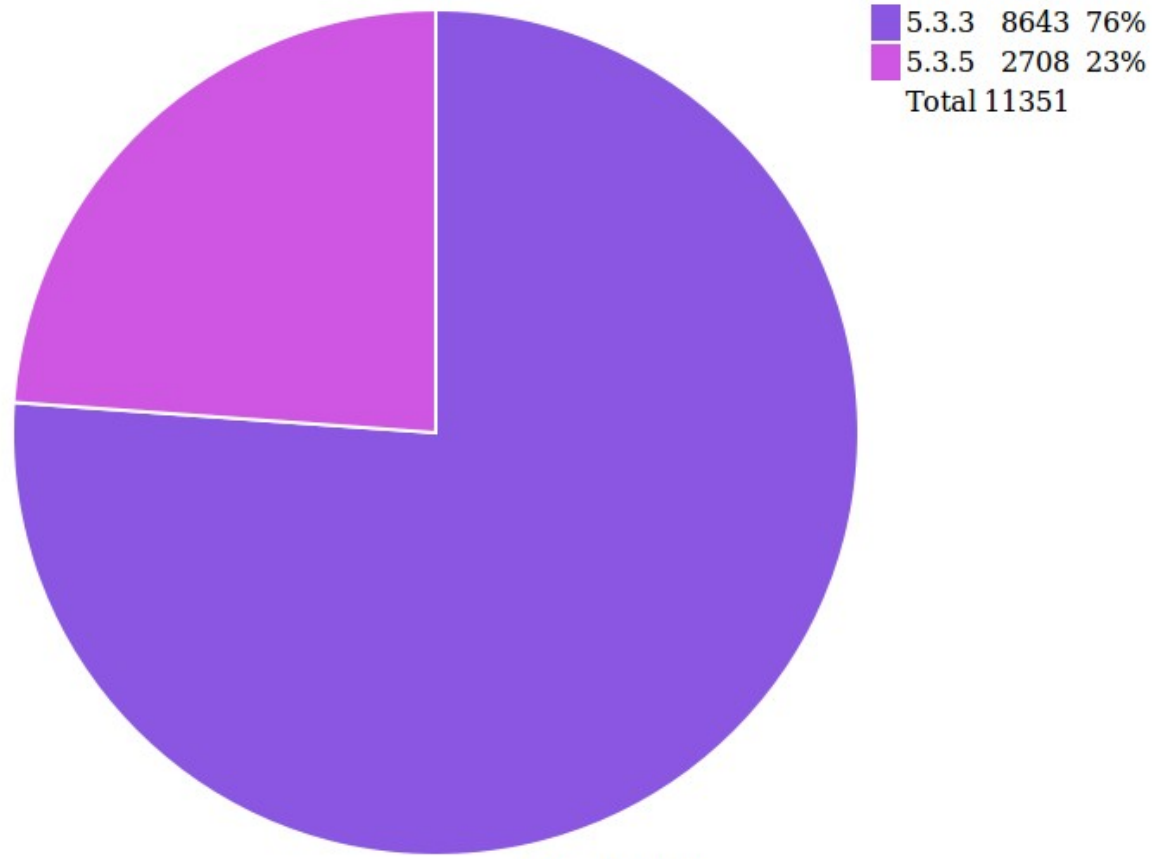
kcureh - program





VASP

Jobs - Version



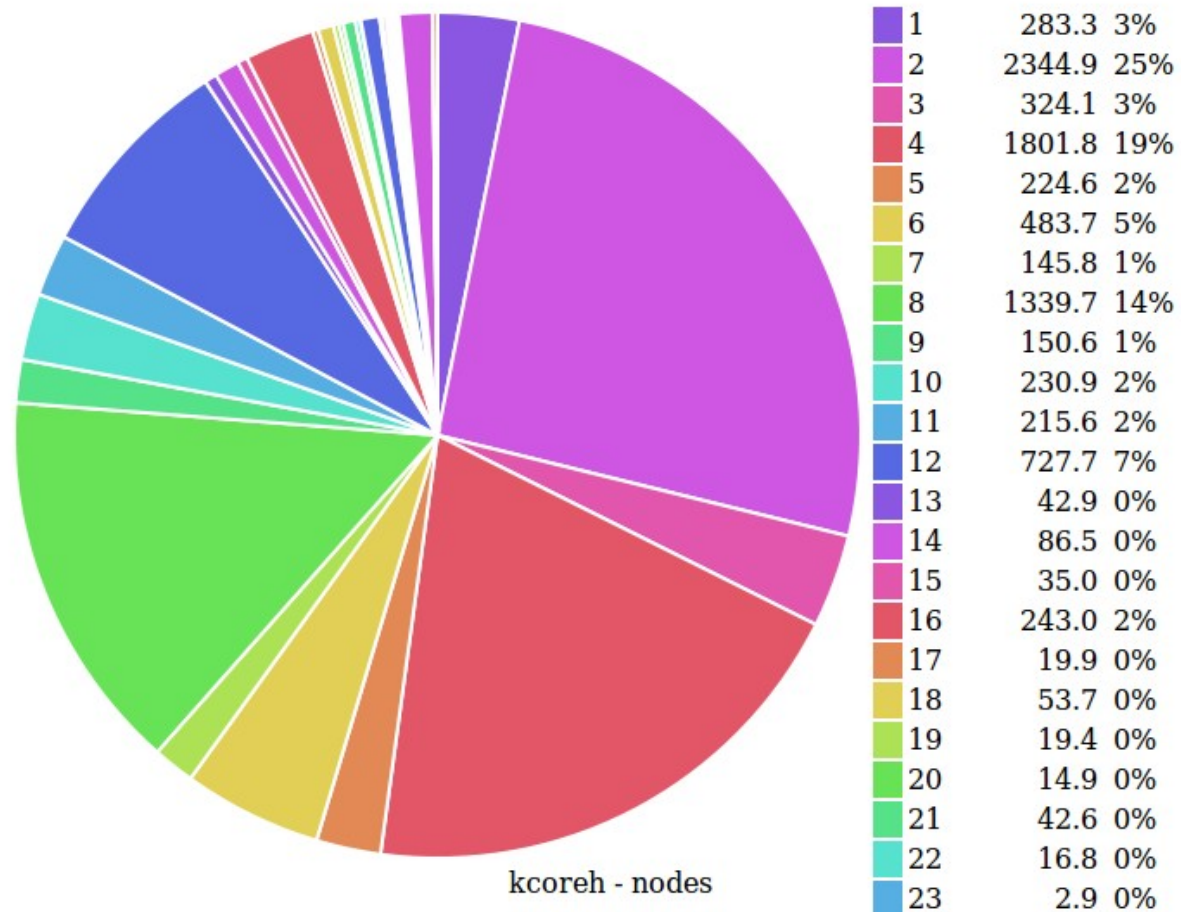
Jobs - Version





VASP

Time - Nodes



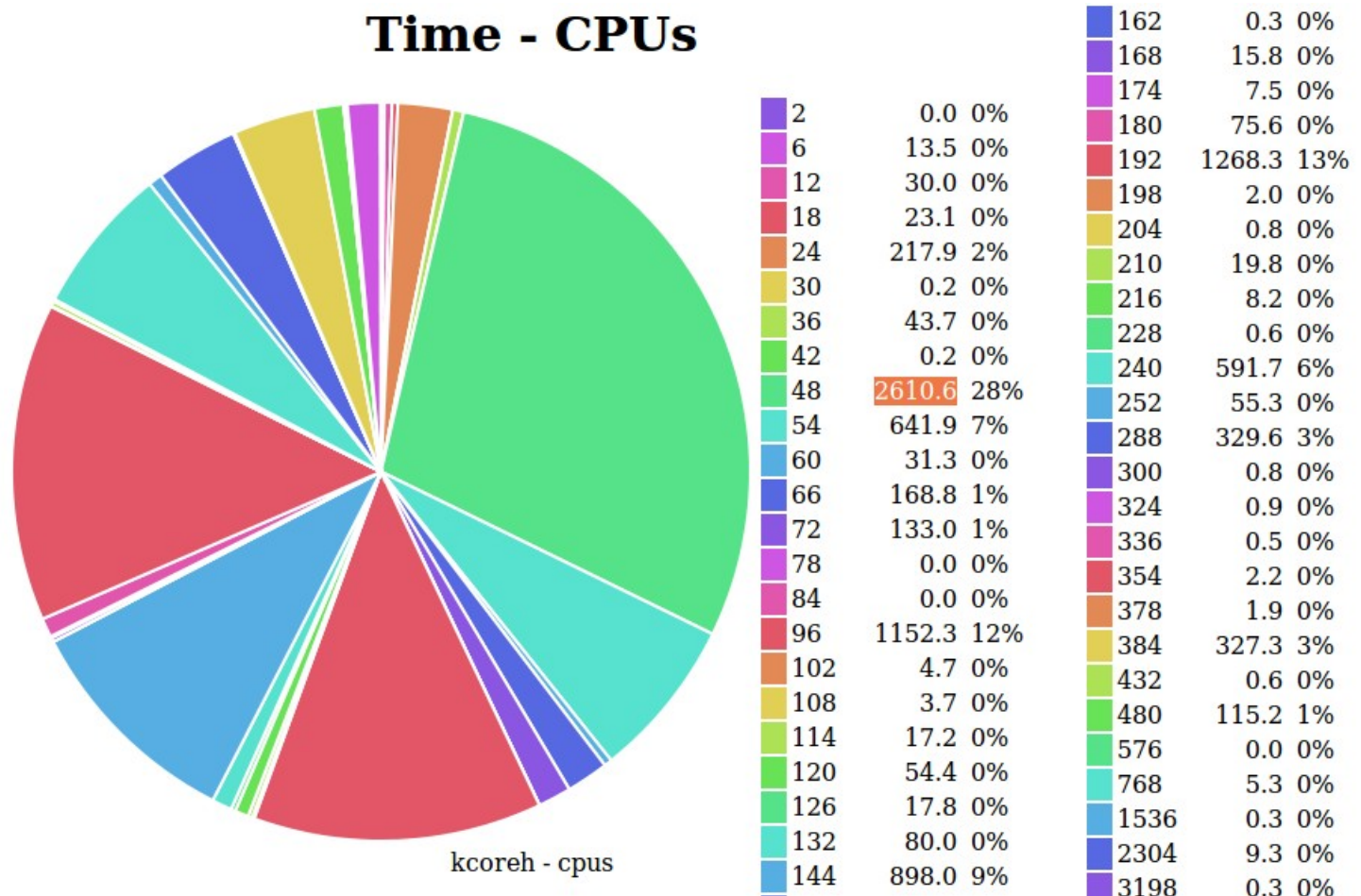
kcoreh - nodes





VASP

Time - CPUs



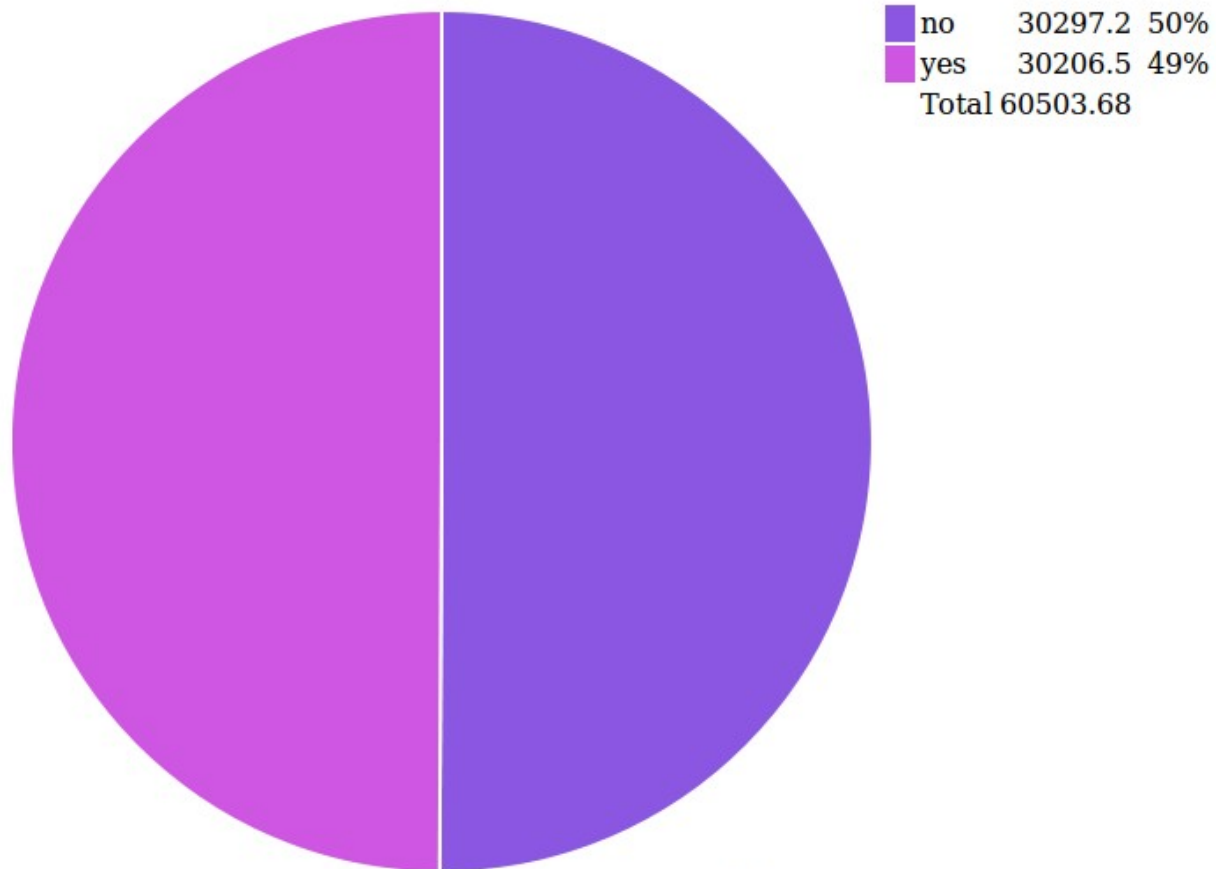
kcoreh - cpus





User Compiled Software

Time - User Compiled



kcureh - User Compiled





Problems?

- User compiled software
 - Do we care?
- When is a job completed?
 - Assume
 - Ask Slurm
- What to save?
- Sample
 - What are we missing?





Live Demo :-)



Thanks for me

- Questions?

