

Improve utilization and load tracking in the scheduler

Presented by
Vincent Guittot

Date
November 3, 2016

Event
Linux Plumber Conference

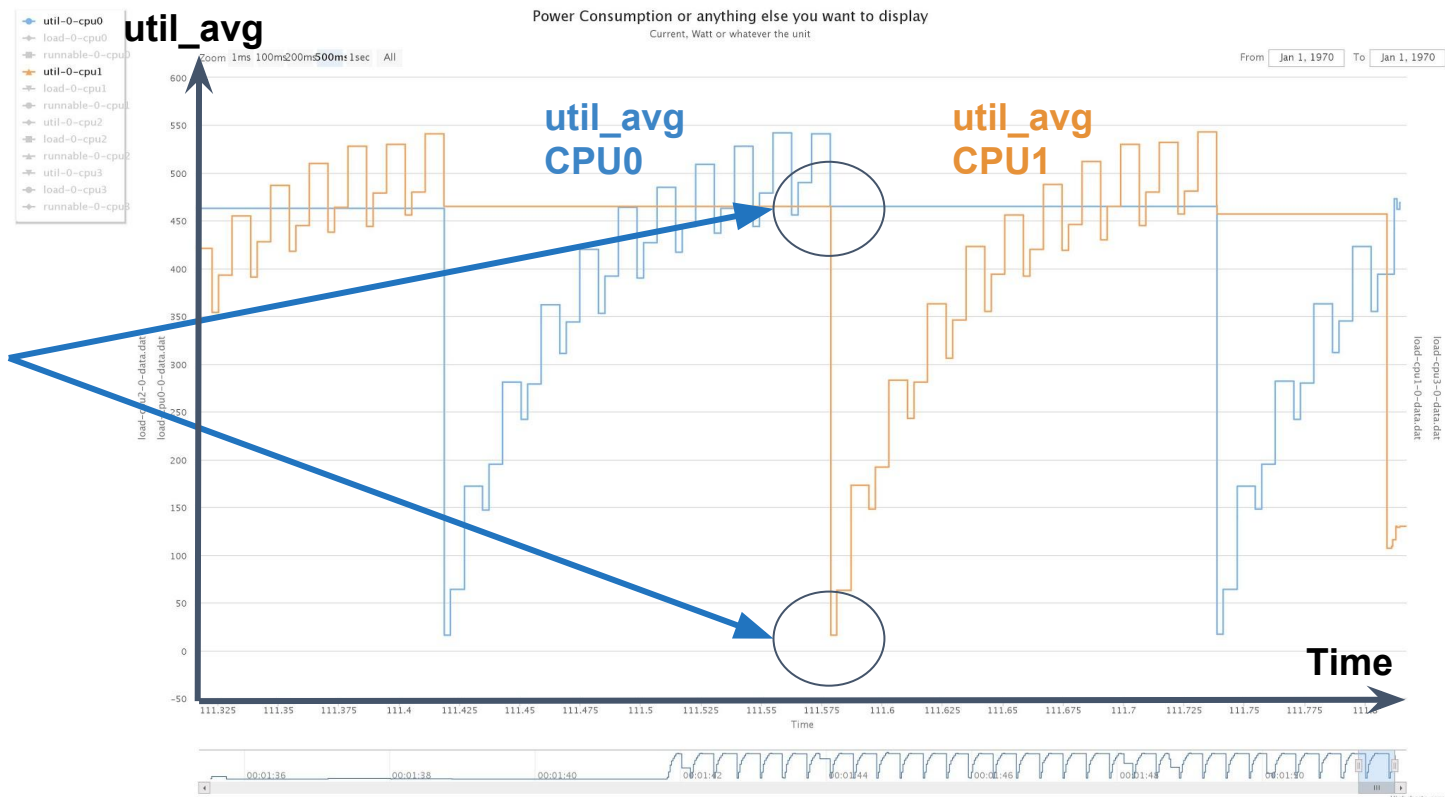
Items

- Propagate tasks movement into task group
- Update idle CPUs
- Scale invariance
- Responsiveness of PELT

Propagate tasks movement into task group

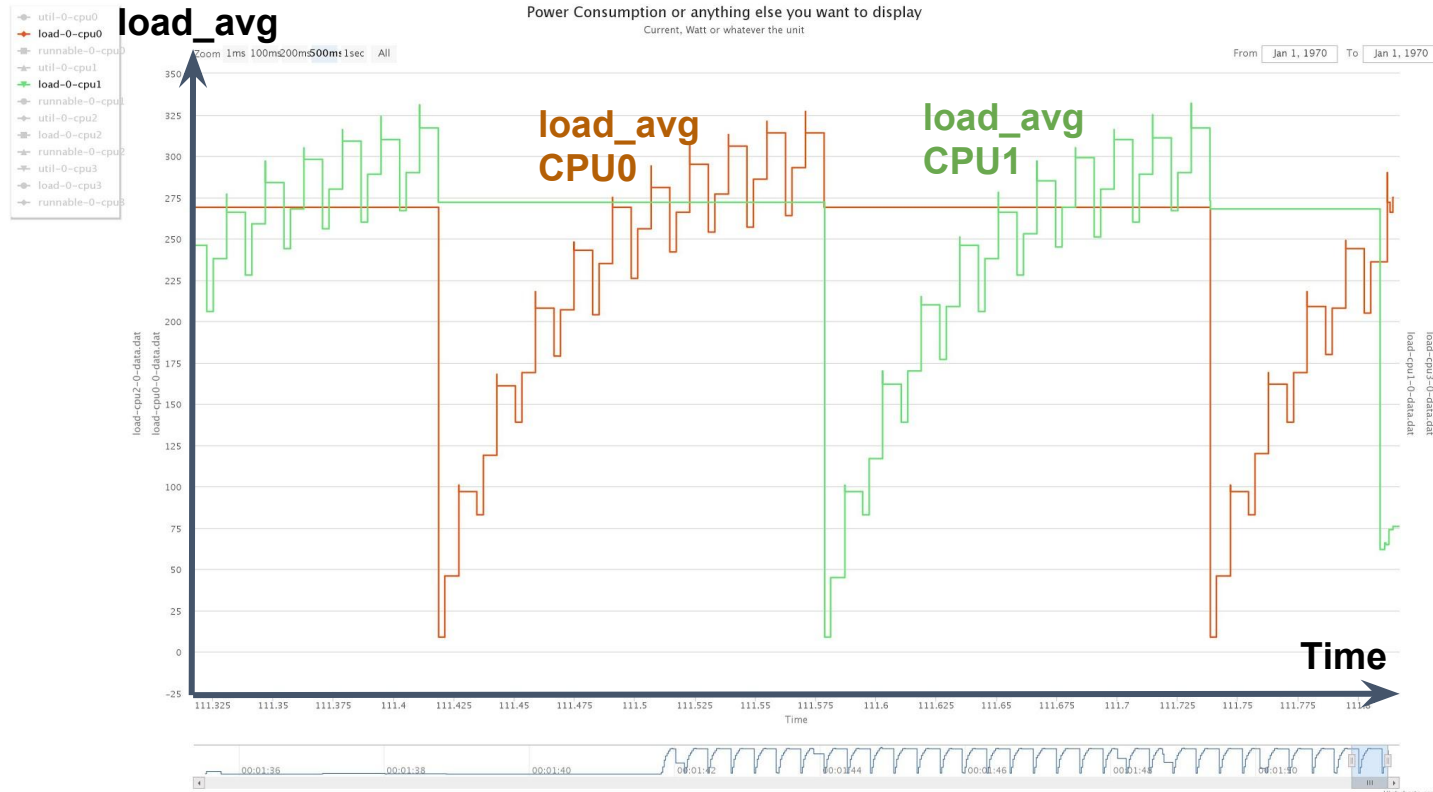
Migrate inside task group

migration is not seen at root level

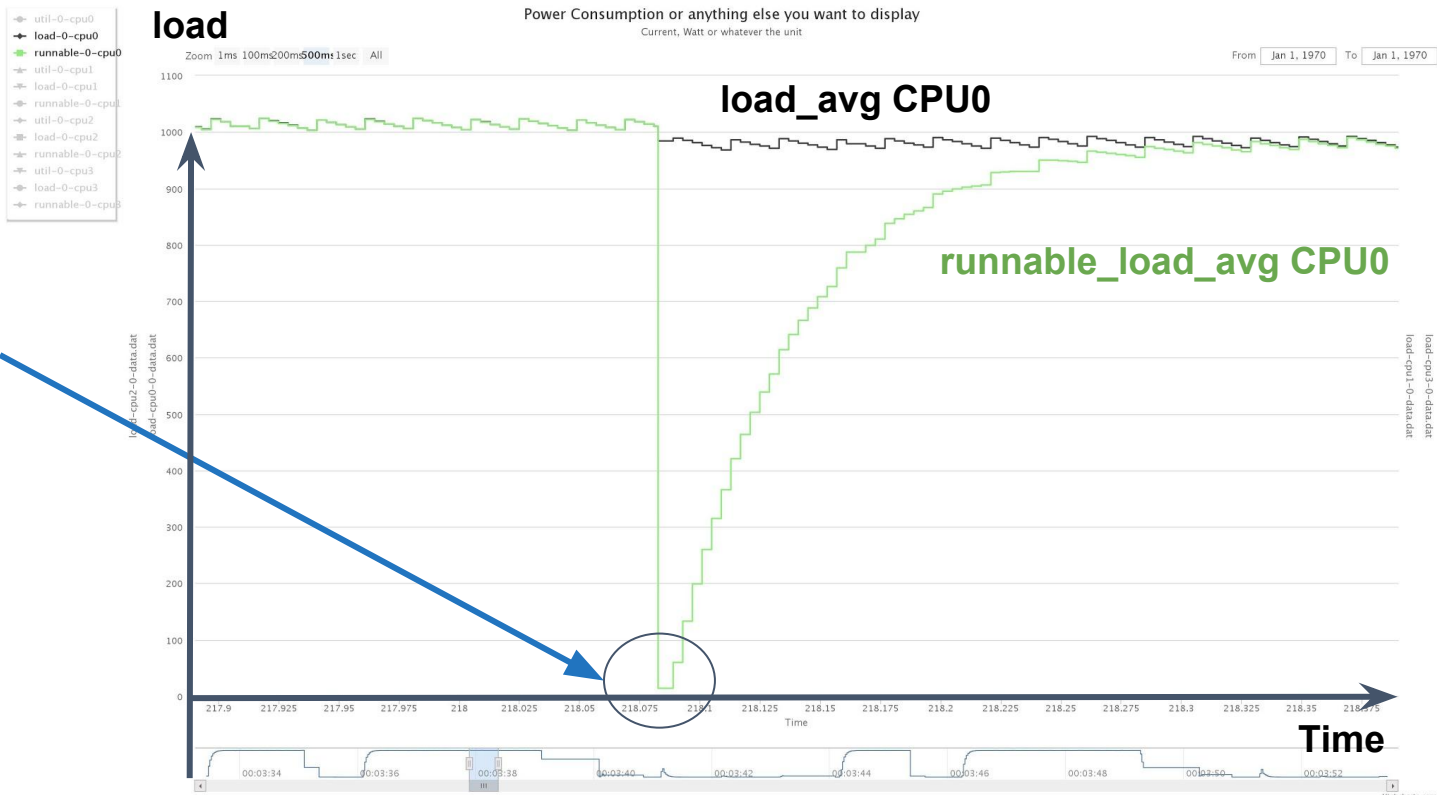


Migrate in task group

Same behavior for load



Move between task groups

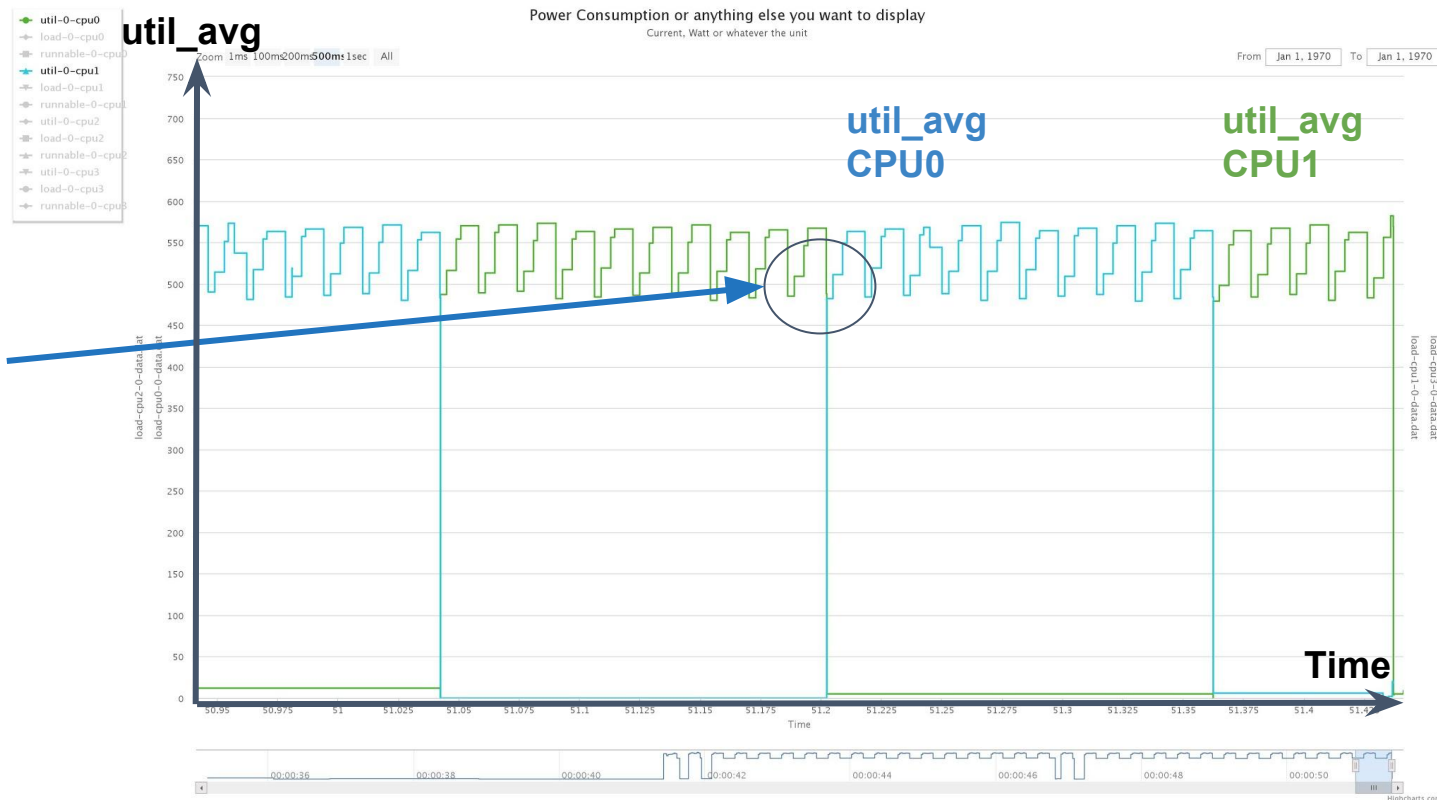


Propagate tasks movement into task group

- Propagate events down to the root domain
 - attach/detach task
 - remove task (asynchronous detach)
 - enqueue/dequeue task (already there)
- Propagate utilization
 - Simply add children changes to parent's utilization
- Propagate load
 - Task group must not consumes more CPU than a task of equal weight
 - Scale task group's load into the range of its shares

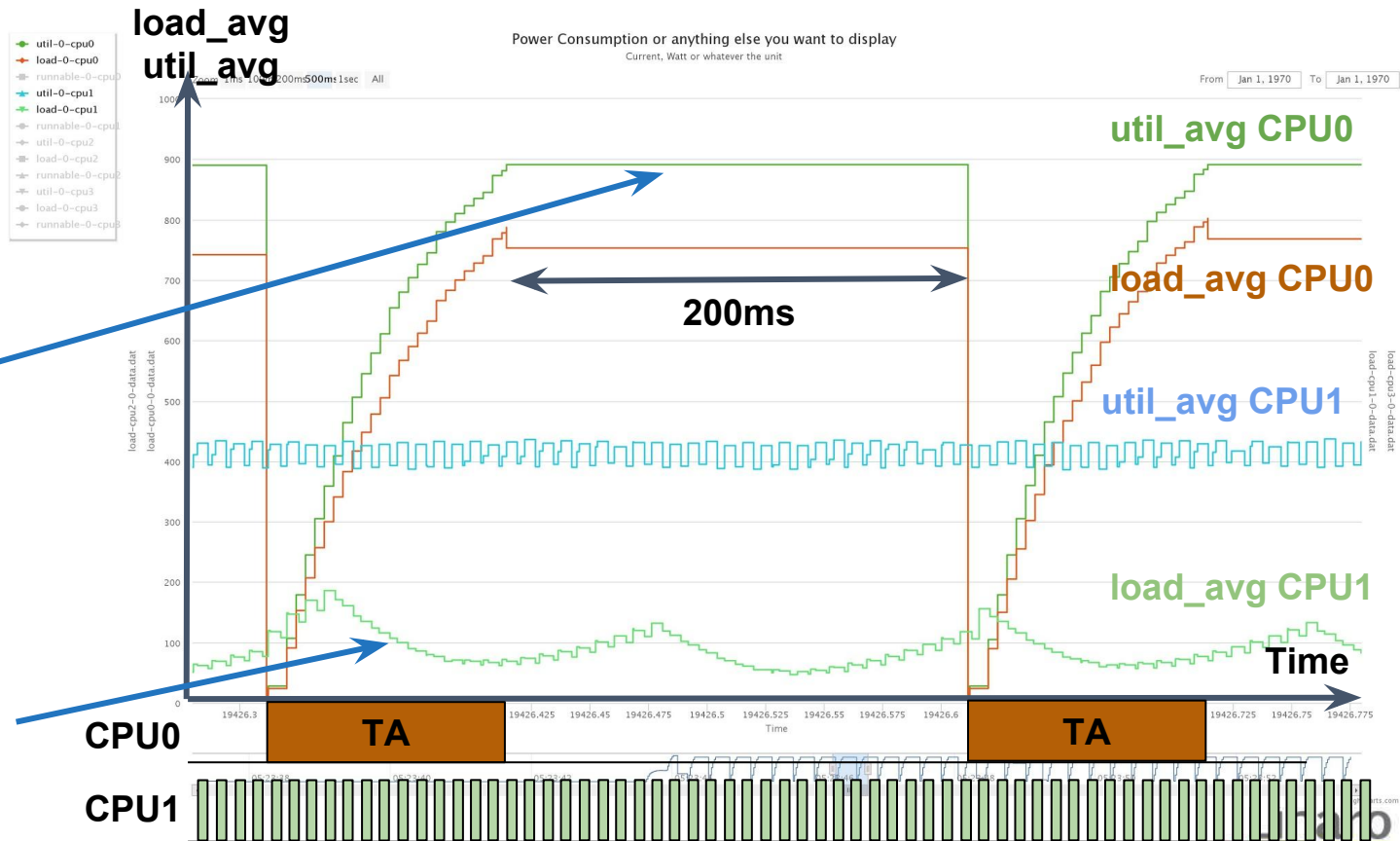
Propagate tasks movement into task group

Utilization migrates
with task



Update idle CPUs

Blocked load/utilization of idle CPUs



load_avg/util_avg stay unchanged until next wake up or idle load balance

Impact the shares of busy CPUs

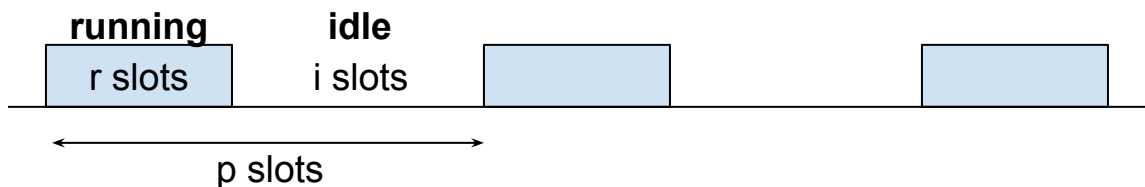


Update idle CPUs

- Ensure periodic update of blocked load
 - On an idle CPU like for periodic idle load balance

Scale invariance

Utilization and load of periodic task

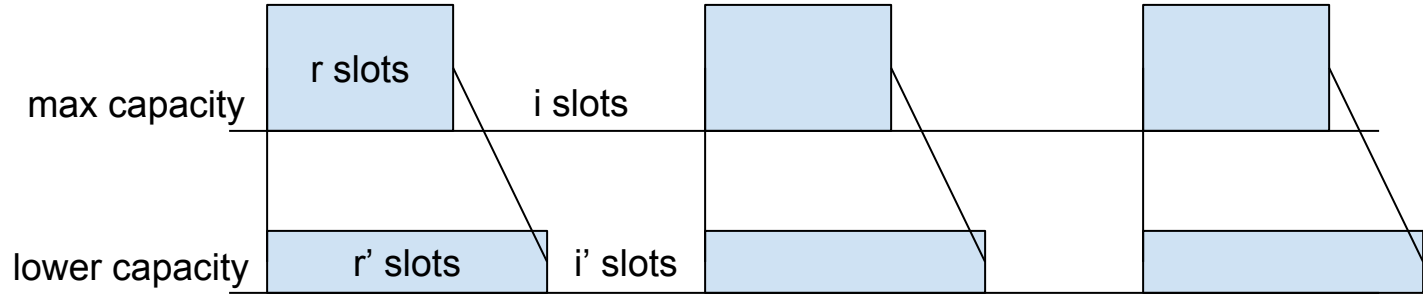


- Utilization range is defined by
 - Number of running slots: r
 - And number of idle slots: i

$$(1-y^r) / (1-y^p) * y^i < \text{Utilization} < (1-y^r) / (1-y^p)$$

- Similar for load

Scale utilization and load



- Invariance applies a scaling factor to the value of time slots

- At max capacity the range is

$$(1-y^r) / (1-y^p) * y^i < \text{Utilization} < (1-y^r) / (1-y^p)$$

- At lower capacity the range is

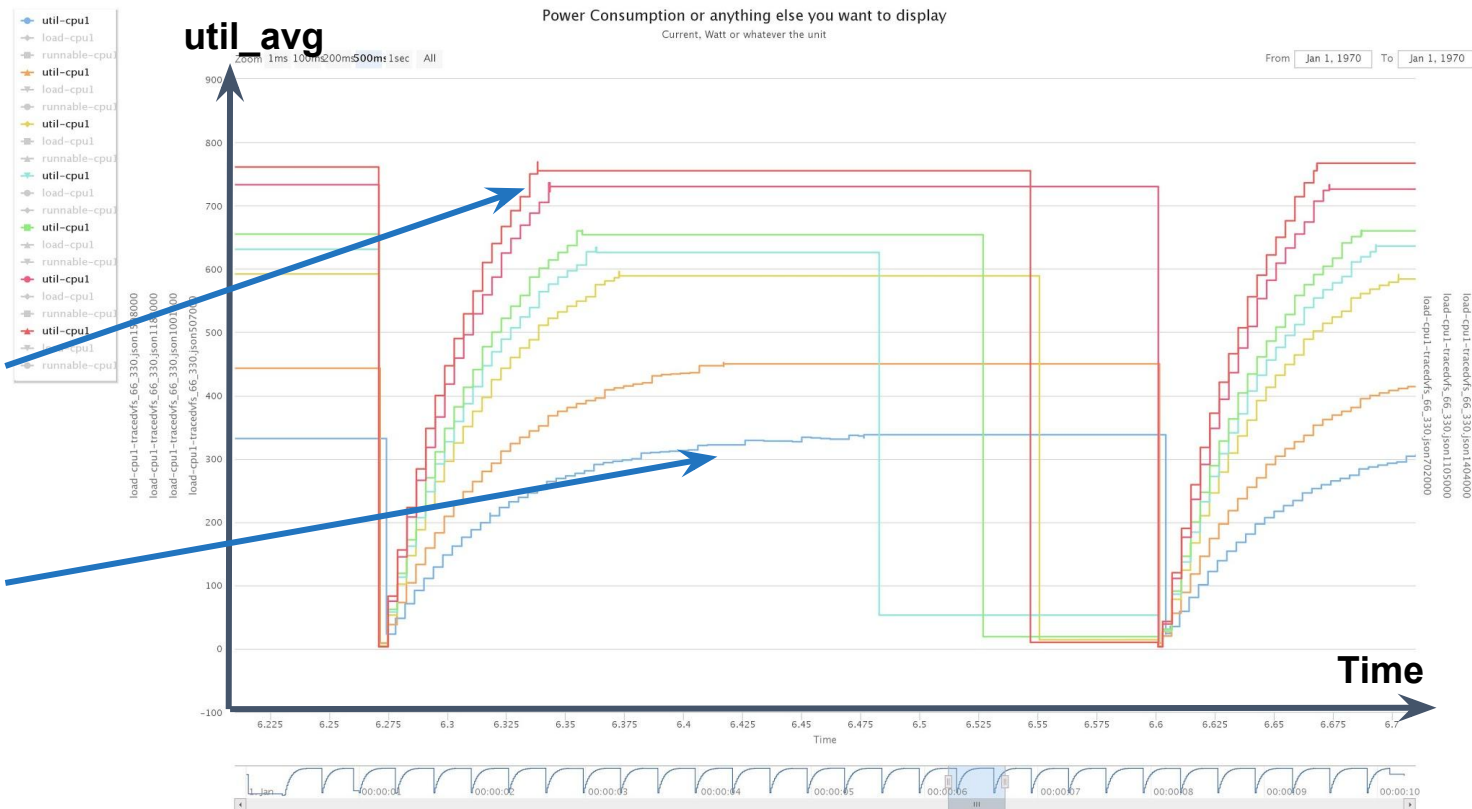
$$C * (1-y^{r'}) / (1-y^p) * y^{i'} < \text{Utilization} < C * (1-y^{r'}) / (1-y^p)$$

with C reflecting the compute capacity ratio

Scale utilization and load

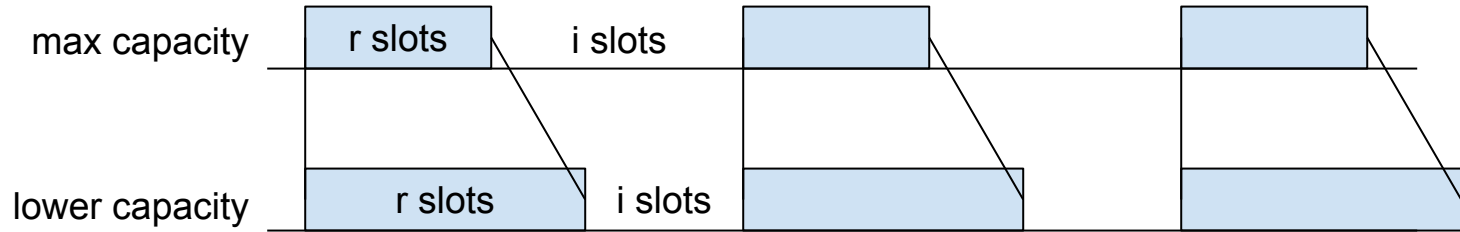
Run 66ms at max capacity

Run 198ms at a 3rd of capacity



Scaled utilization for a task running 66ms every 330ms

Scale time



- Scale duration of time slot instead of scaling utilization of time slot

Impact on load/utilization

- Utilization/Load no more capped by current capacity
 - they can reach max value but need more time
- Utilization/Load of a task is fully invariant
 - Same min/max value at wake up/sleep whatever the compute capacity

Capacity	Max	Half (mainline)	Half (new)	A third (mainline)	A third (new)
Util(%)					
972 (95%)	138ms	N/A	276ms	N/A	414ms
486 (47.5%)	30ms	138ms	60ms	N/A	90ms
256 (25%)	13ms	32ms	26ms	64ms	39ms

Responsiveness of PELT

Responsiveness of PELT

- Some figures
 - Time to reach 95%: 140ms
 - Time to reach 50% : 31ms
- Current decisions are proportional to utilization value
 - High Utilization value selects High OPP
 - Low Utilization value selects Low OPP
- Use other properties of utilization signal ?
 - Use variation
 - An increase of 100 in one running time means that it has run for at least 5ms
 - Apply some kind of PID algorithm
 - take into account variation of utilization

Other alternatives to improve responsiveness

- Use other scheduling class
 - deadline scheduler
- Run at f_{max} for some tasks
 - Boost utilization of some tasks and group of tasks
- Keep track of “peak per-entity utilization”
 - Track last max utilization value
- Use other load tracking mechanism

