

# Rework load\_balance

Vincent Guittot

20 May 2019



# Agenda

- Introduction
- UCs
- Statistics
- Policy

# Introduction

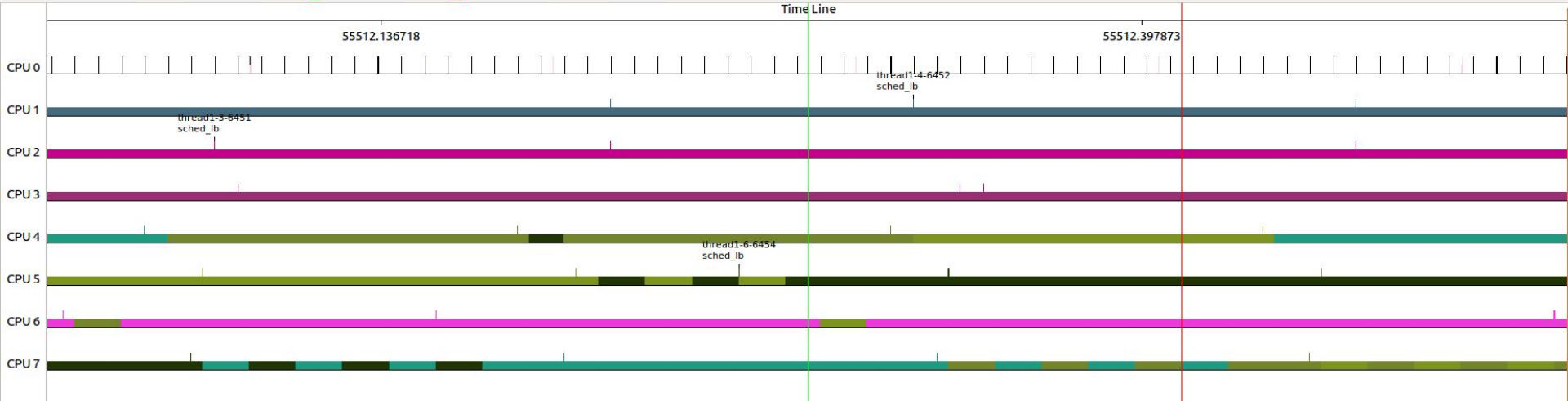
- Some UCs not correctly handle yet
- Others seem to be randomly supported
  
- Simplify calculate\_imbalance
  - Remove old/meaningless heuristics
  - Remove weird imbalance calculation
    - Don't set a value that "should" be enough to fix the imbalance
  
- Replace UC's specific classification of group by generic tag
  
- Metrics have evolved
  - Runnable\_load\_avg
  - Util\_avg
  - Load\_avg

# Some UCs

# 1 task per CPU

- 1 task per CPU with asymmetries
  - big Little

Pointer: 55512.263955 Cursor: 55512.283645 Marker 55512.283709 Marker 55512.411817 A,B Delta: 0.128108



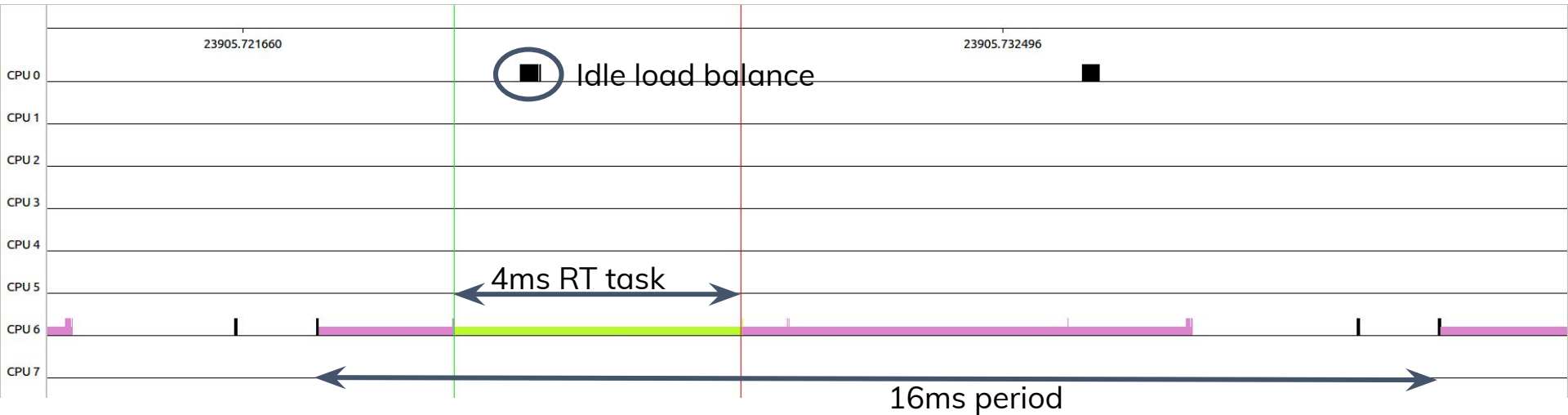
- RT/DL/IRQ stolen time
- Thermal pressure ?

# fix\_small\_imbalance

- Depends of the weight of the tasks
- Depends of the initial state
  - Tasks placement
- Asymmetries screw up the default load balance mechanism
  - Compare load\_avg

# Class preemption

- CFS task preempted by RT but idle/newidle doesn't pull waiting task
  - Only takes into account cfs.h\_nr\_running



# Running tasks

- Take into account tasks from other classes
  - Minimize resources contention
- Take into account sched\_idle task
  - Improve scheduling latency



# Ensure fairness

- Unbalanced system
  - 9 tasks on 8 CPUs
- When/How do we periodically migrate tasks to ensure fairness ?



# Ensure fairness

- When does a task migrate ?
  - At every load\_balance
  - with fix\_small\_imbalance
- Result looks a bit random
  - Use a more stable mechanism

# Some UCs

- And probably lot of others

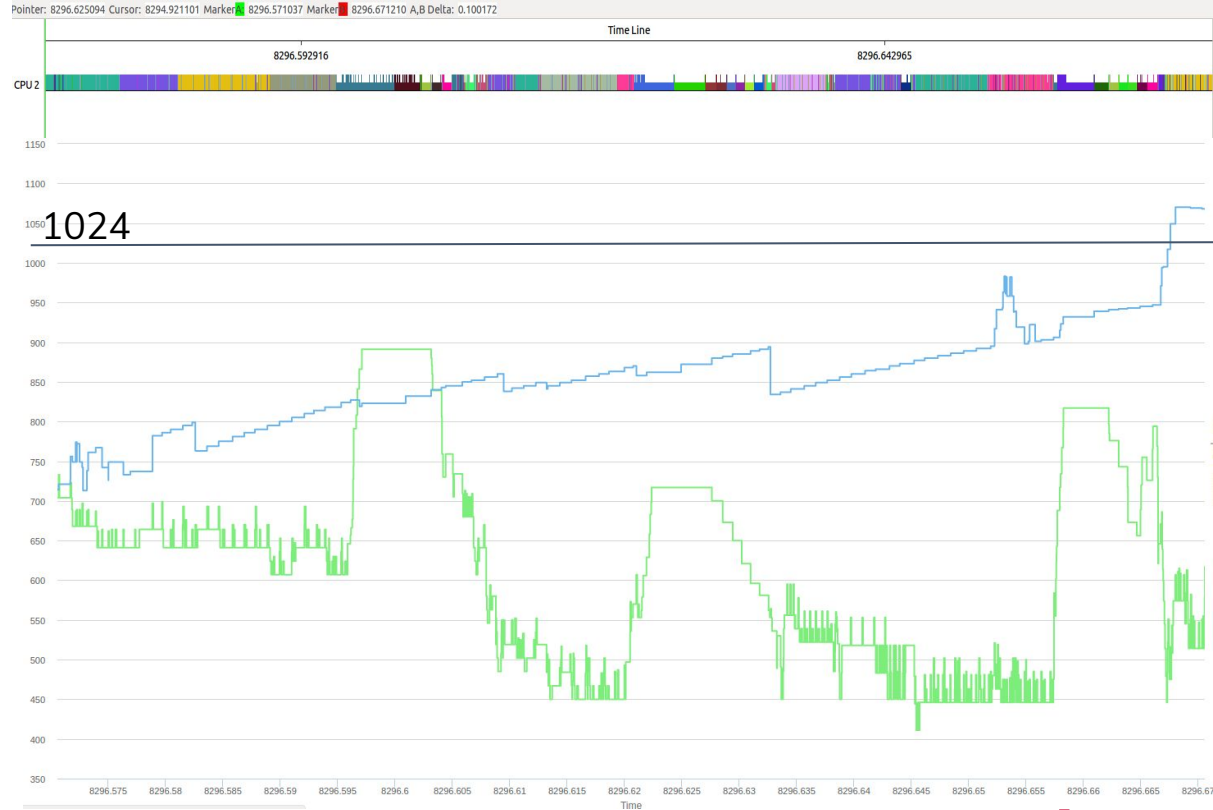
# Statistics

# Statistics

- running tasks
  - only nr\_h\_running is tracked
  - What about nr\_running ?
  - What about CPUs with only sched\_idle ?
- Utilization:
  - Util\_avg is used to estimate spare capacity
  - But take care of jump due to migration
  - Util\_est not always mitigate this effect
- Load :
  - Runnable is currently used
  - Lot of variance
  - Track both load and runnable\_load ? similarly to find\_idlest\_group
- Capacity
  - Average available capacity: Max - IRQ/DL/RT and may be thermal

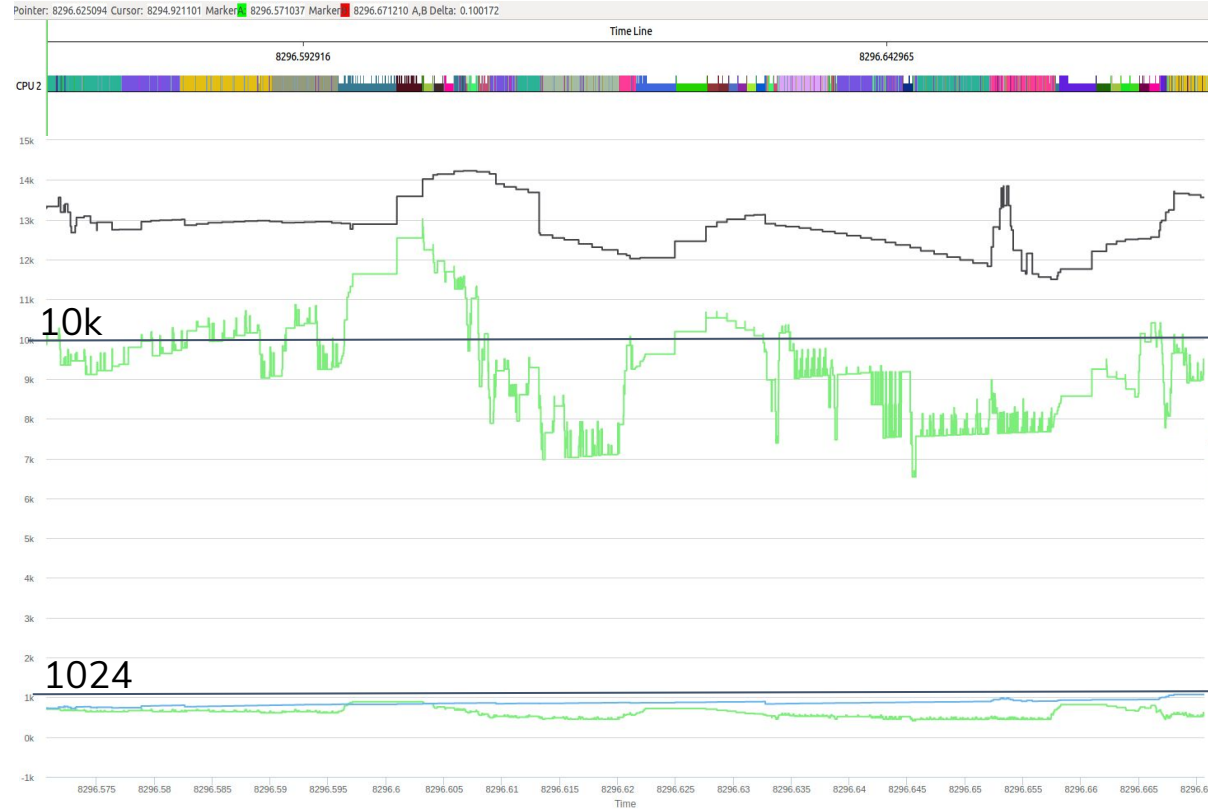
# utilization

- Hackbench
  - CPUs are fully used
  - Lot of migration
- util\_avg or util\_est
  - Not always meaningful
- Take into runnable time ?
  - waiting time



# runnable load

- runnable load quite high
- Lot of waiting time



# Runnable vs load

- Runnable\_load vs load
  - Runnable\_load gives current running state
  - can use nr\_running and utilization instead ?
- Remove weight from runnable
  - Track waiting time of tasks



# Statistics

- Don't bias statistics like with prefer sibling
  - Tag the group instead
- load\_per\_task
  - Doesn't mean anything
  - Used when we don't know how much to move
  - Should be replaced by an explicit force migration
- sum\_weighted\_cpuload
  - Is the same as group\_load since 4.19 and the disable of LB\_BIAS
  - And remove cpu\_load ...

# Statistics

- nr\_h\_running
- nr\_running
- sched\_idle state
- Utilization :  $\max(\text{util\_ag}, \text{util\_est})$
- Runnable\_load\_avg
- Load\_avg
- capacity

# Policy

# Load balance sequence

- Load\_balance
  - should\_we\_balance: select 1st idle or 1st CPU of the group
  - **find\_busiest\_group**(&env);
  - find\_busiest\_queue(&env, group);
  - Move waiting tasks or running task
  - Update balance interval
  
- Most of the work is done in
  - find\_busiest\_group
  - And find\_busiest\_queue to a lesser extend

# Find busiest group

- `find_busiest_group`
  - **Update statistics and classify groups**
  - A list of conditions to decide if we are balanced
  - Calculate imbalance

# Classify group

- Default state for the group
  - other
  - overloaded
- Plus some special cases
  - Imbalance
  - Misfit task
- Plus some short cut for asym packing
- Simplify the algorithm
  - Remove UC dedicated state

# Group state

- Set a state for the group
  - Has capacity
  - Fully Busy
  - Overloaded
- Plus special cases
  - Fill an empty cpu
  - Move any task but at least one
  - Move a minimum load/util
  - Move a dedicated task

# Policy

- Normal case: “overloaded”
  - we compute the load to move
- Other specific UCs where we want to move task
  - Fill an empty cpu
  - Unblocked a pinned related imbalance by moving any task
  - Move big task on big CPU
  - A dedicated task
- Imbalance can be a load or an utilization or a dedicated task



# Thank you

Join Linaro to accelerate deployment of your  
Arm-based solutions through collaboration

[contact@linaro.org](mailto:contact@linaro.org)





