

# Energy-aware Scheduling

국민대학교 임베디드 연구실  
경 주 현

# Outline

- **Summarize two documents**
- **“Energy-aware Scheduling”**
- **“EAS Overview and Integration Guide”**

# Motivation

- **Hardware topologies are becoming more varied**
- **accommodating different power/performance budgets:**
  - SMP, multi-cluster SMP, ARM big.LITTLE technology.
  - Per core/per cluster DVFS
- **Linux power management frameworks are uncoordinated**
  - Hard to tune for different topologies.

# Power Saving Scheduling

- **Measure all CPUs' load**
- **Choose a CPU from available CPUs**
- **Decide the idle CPU**
- **Migrate tasks from the idle CPU to the chosen CPU**

# State of the art document for EAS.

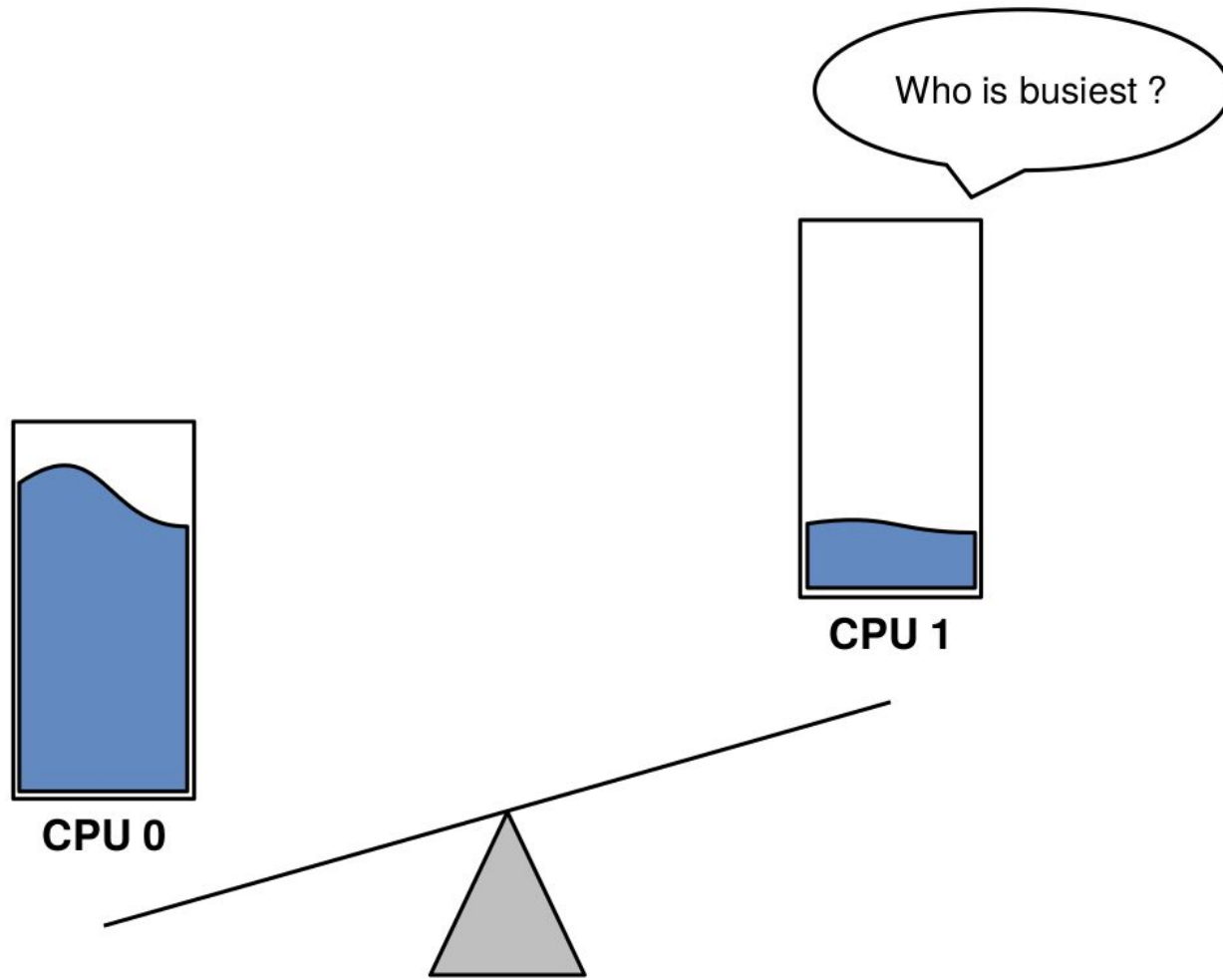
- **"Energy Aware Scheduling", Byungchul Park, LG Electronic**
- **"EAS Overview and Integration Guide", ARM**



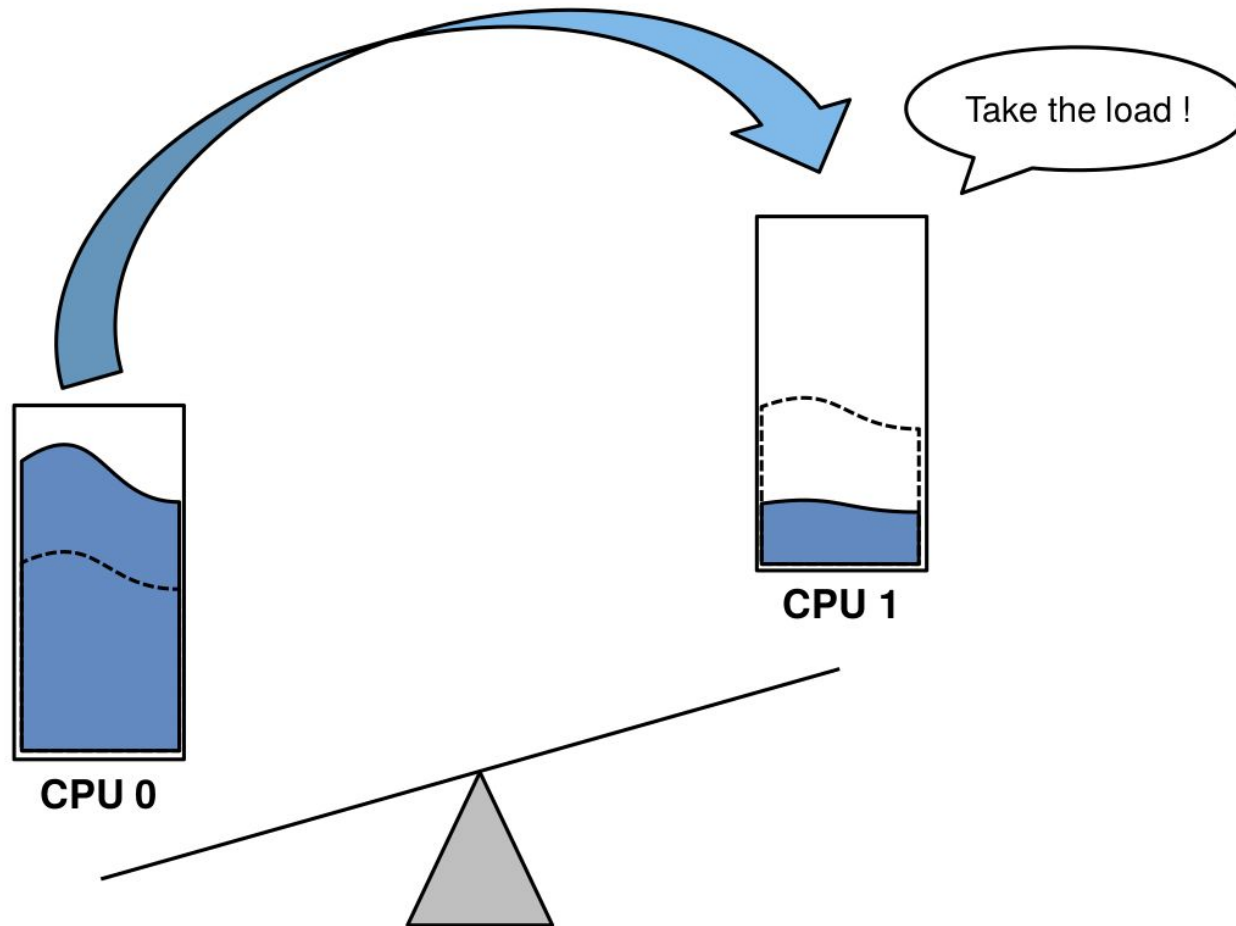
**"Energy Aware Scheduling",  
Byungchul Park,  
LG Electronic**

**reference link: [Linux Foundation Events](#)**

# Load Balancing - migration



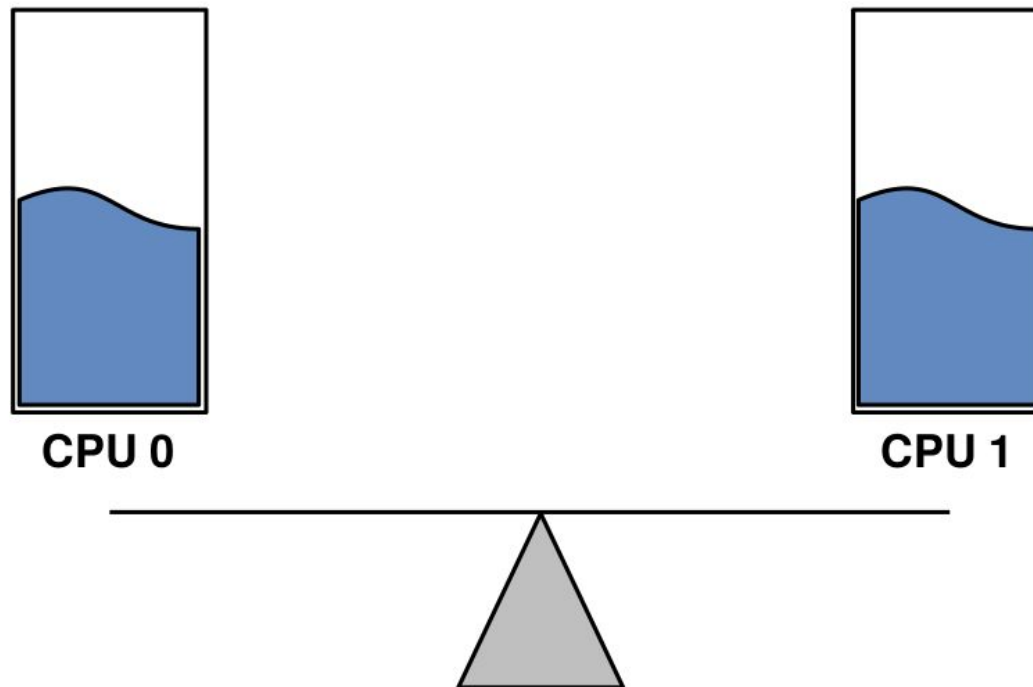
# Load Balancing - migration



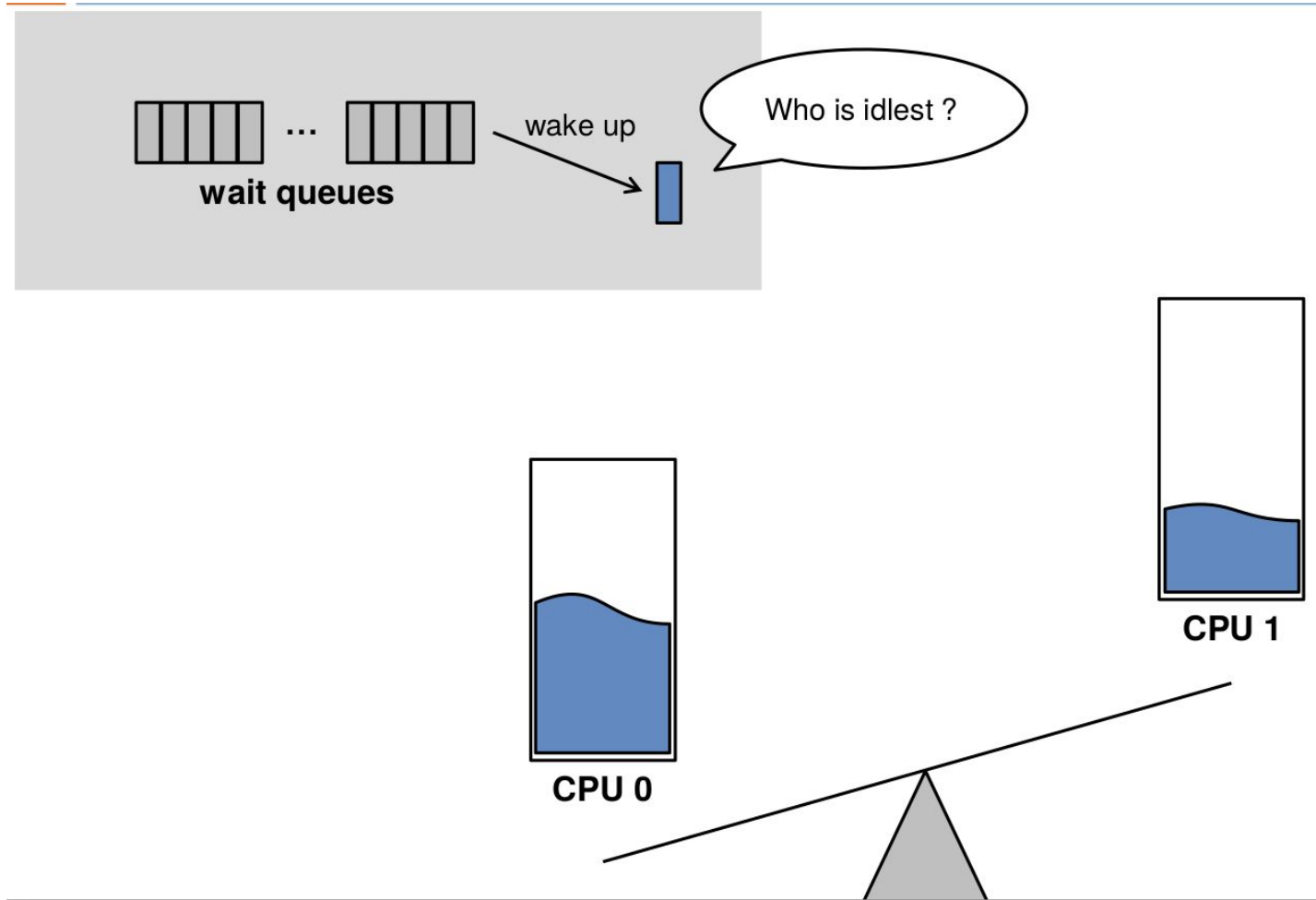


# Load Balancing - migration

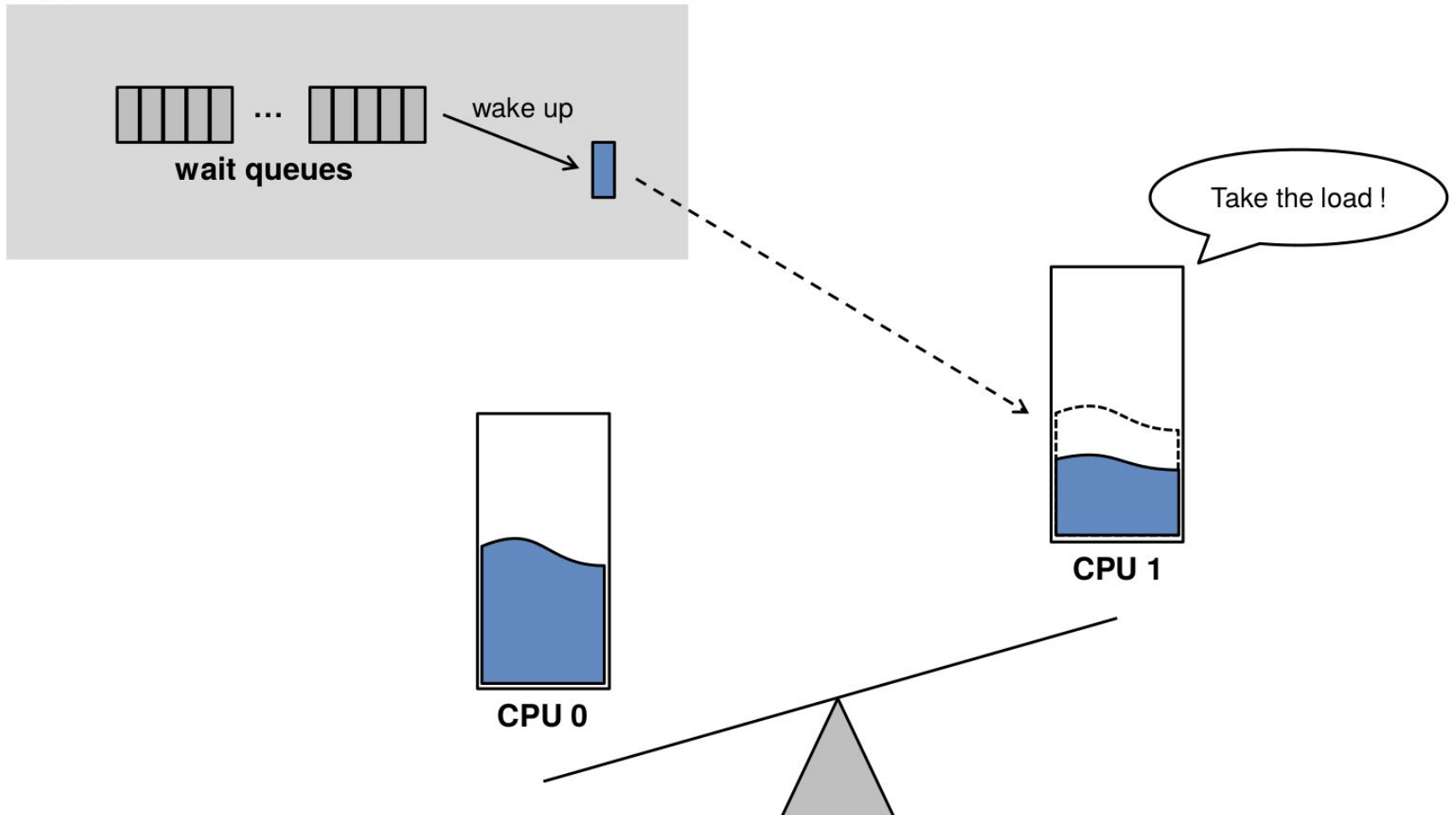
BALANCED !!!



# Load Balancing - wake-up



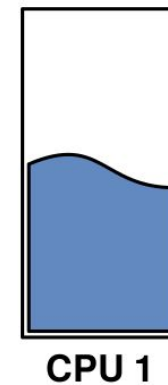
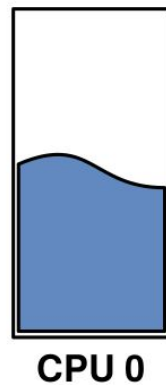
# Load Balancing - wake-up



# Load Balancing - wake-up



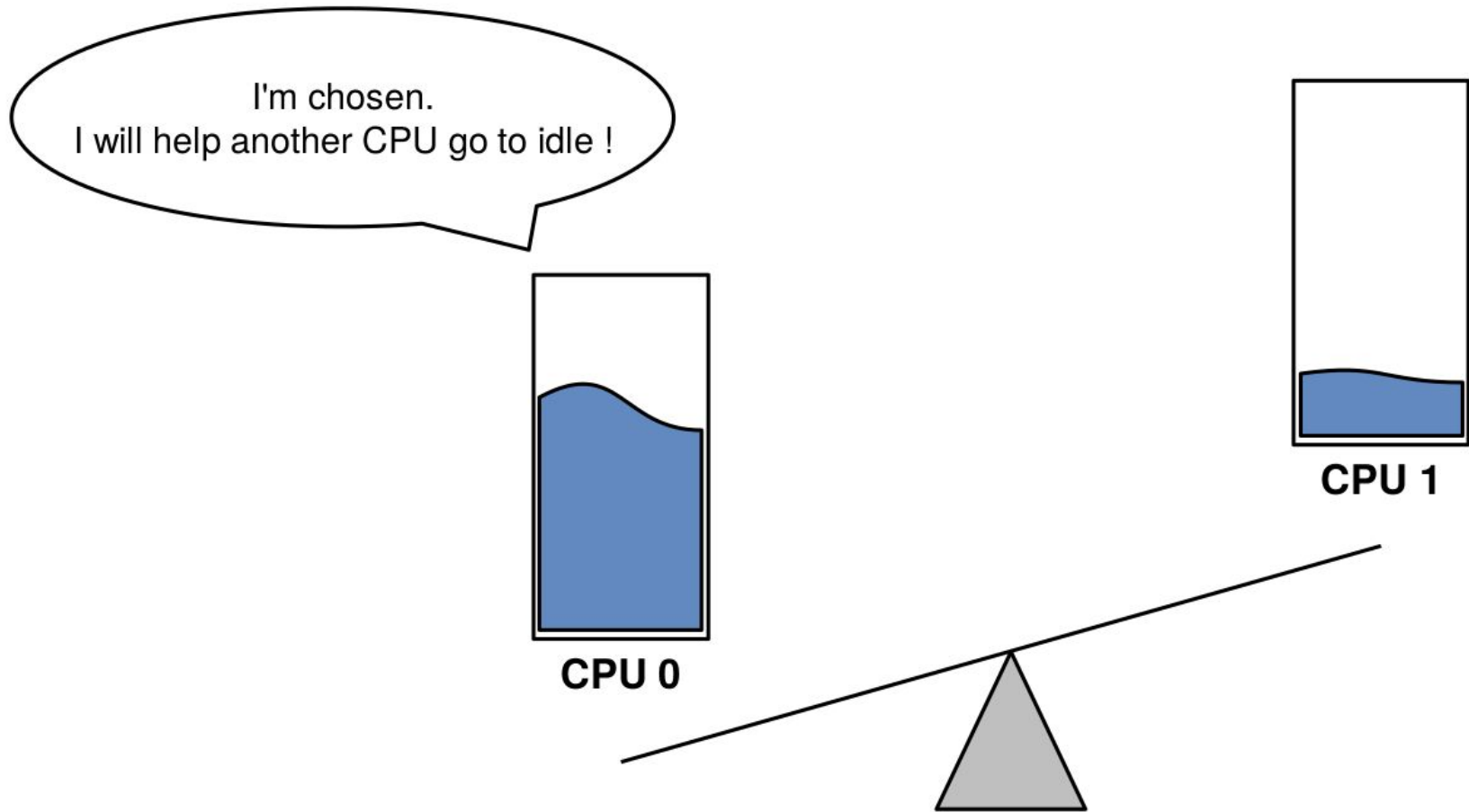
BALANCED !!!



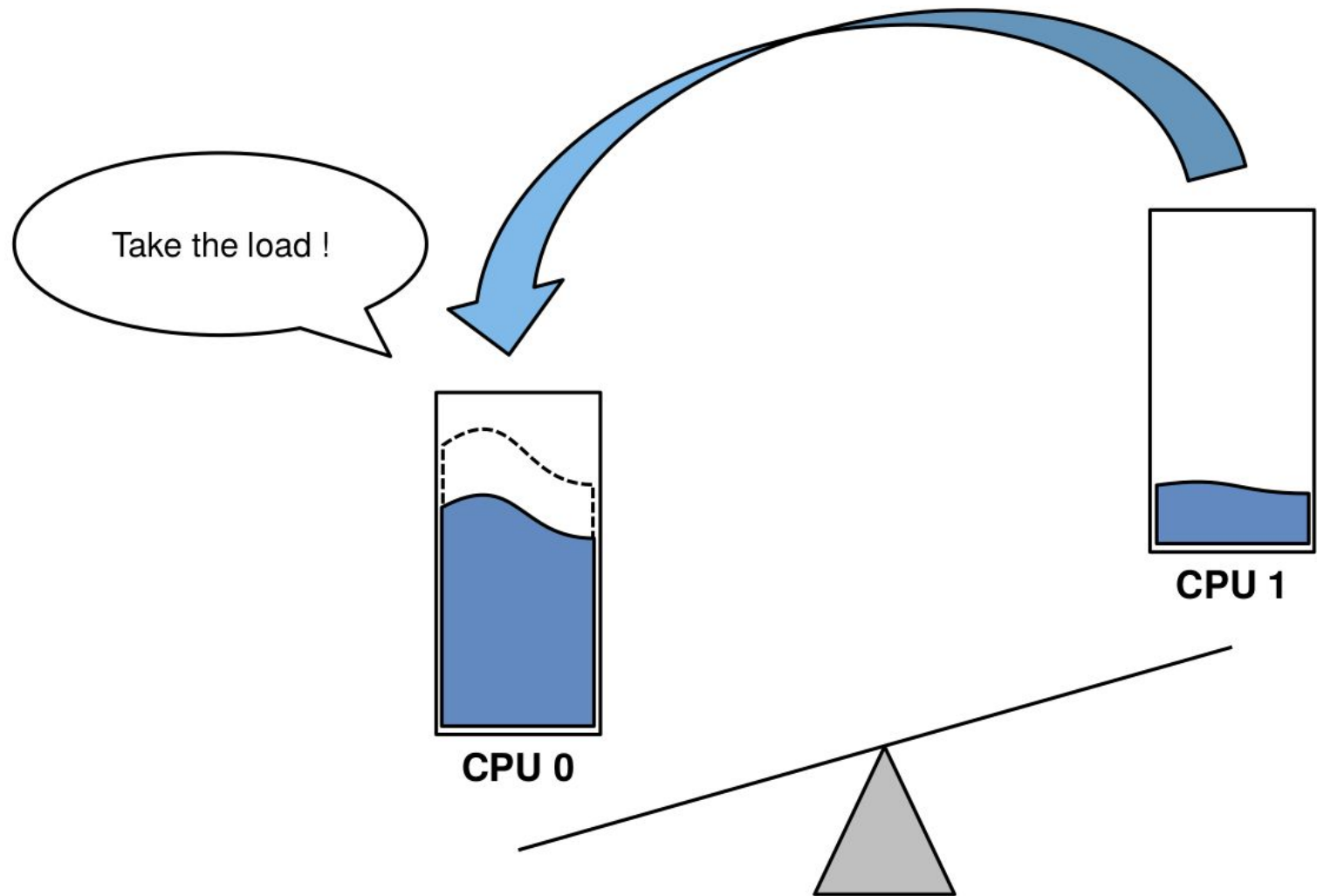
# Problem

- **Only focused on performance**
- **Spreading approach**
  - Guarantee maximum performance and minimum latency

# Power Saved Load Balancing - migration

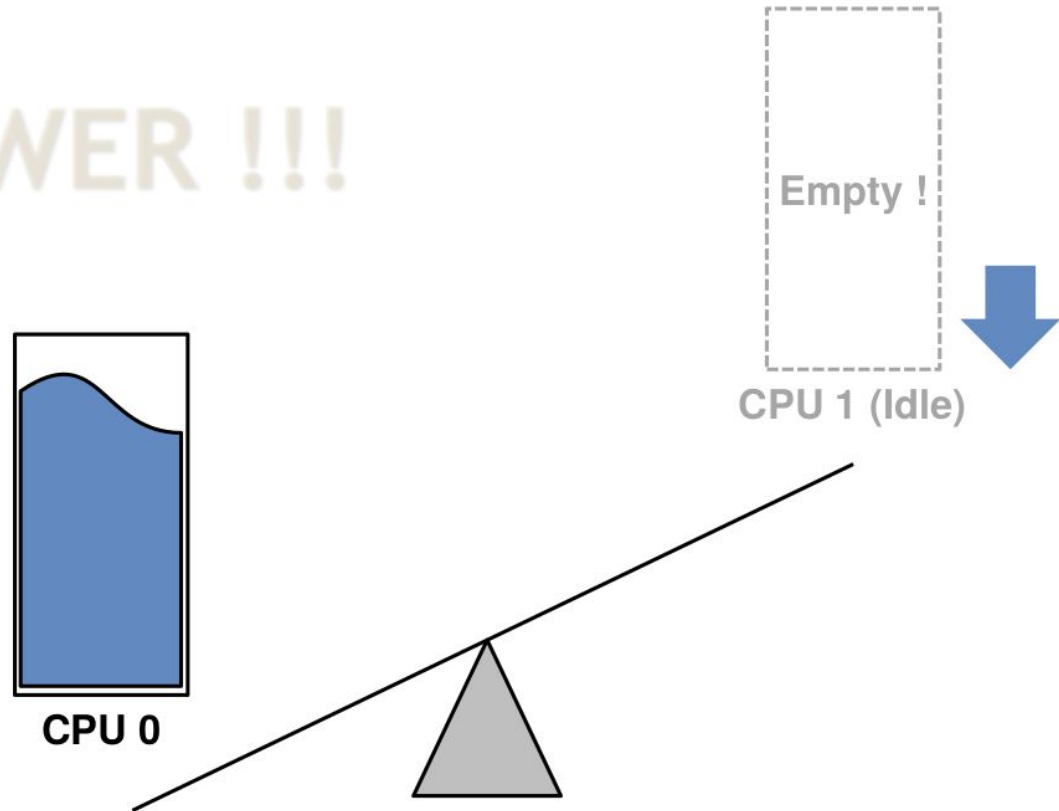


# Power Saved Load Balancing - migration



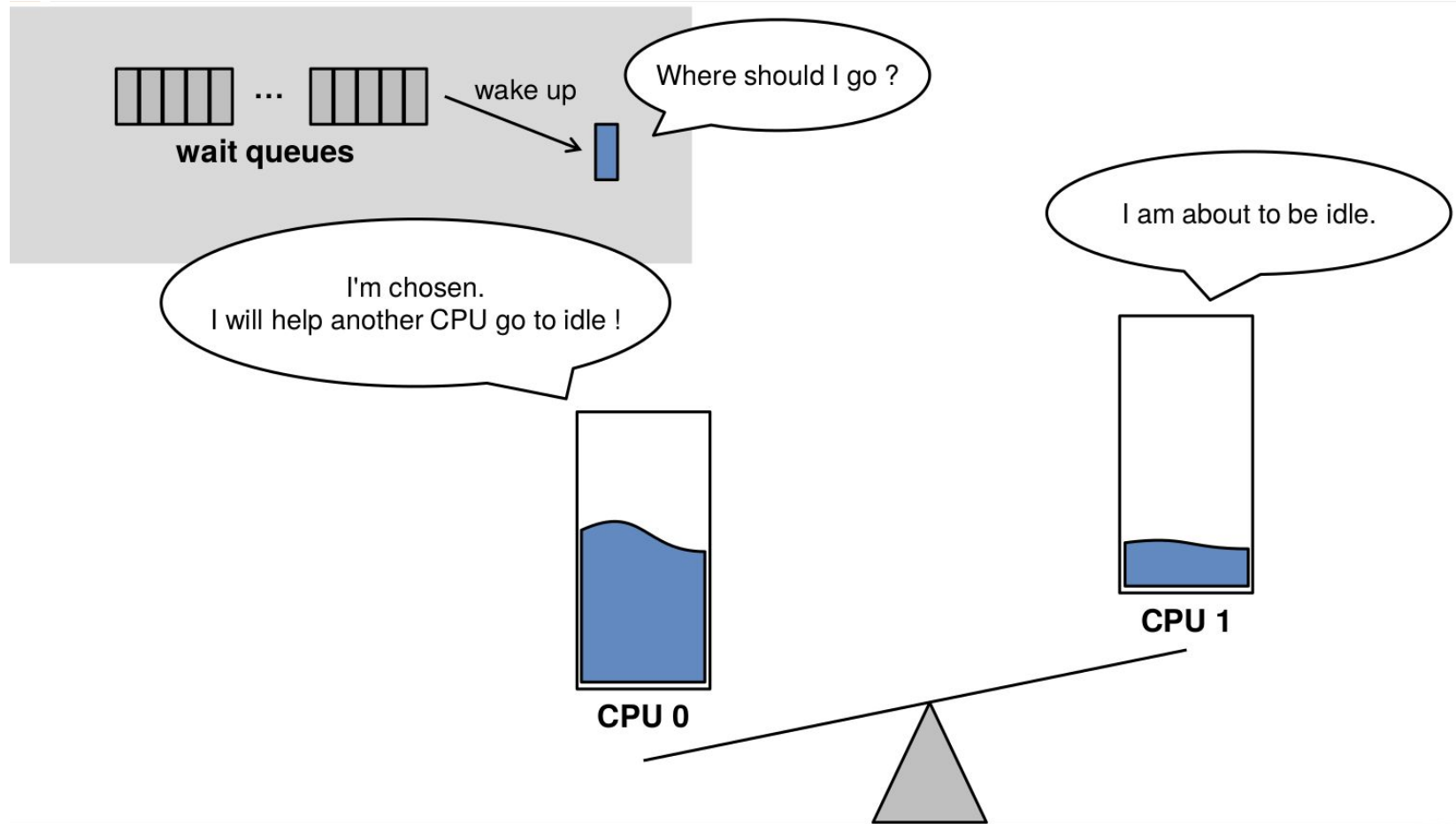
# Power Saved Load Balancing - migration

SAVE POWER !!!

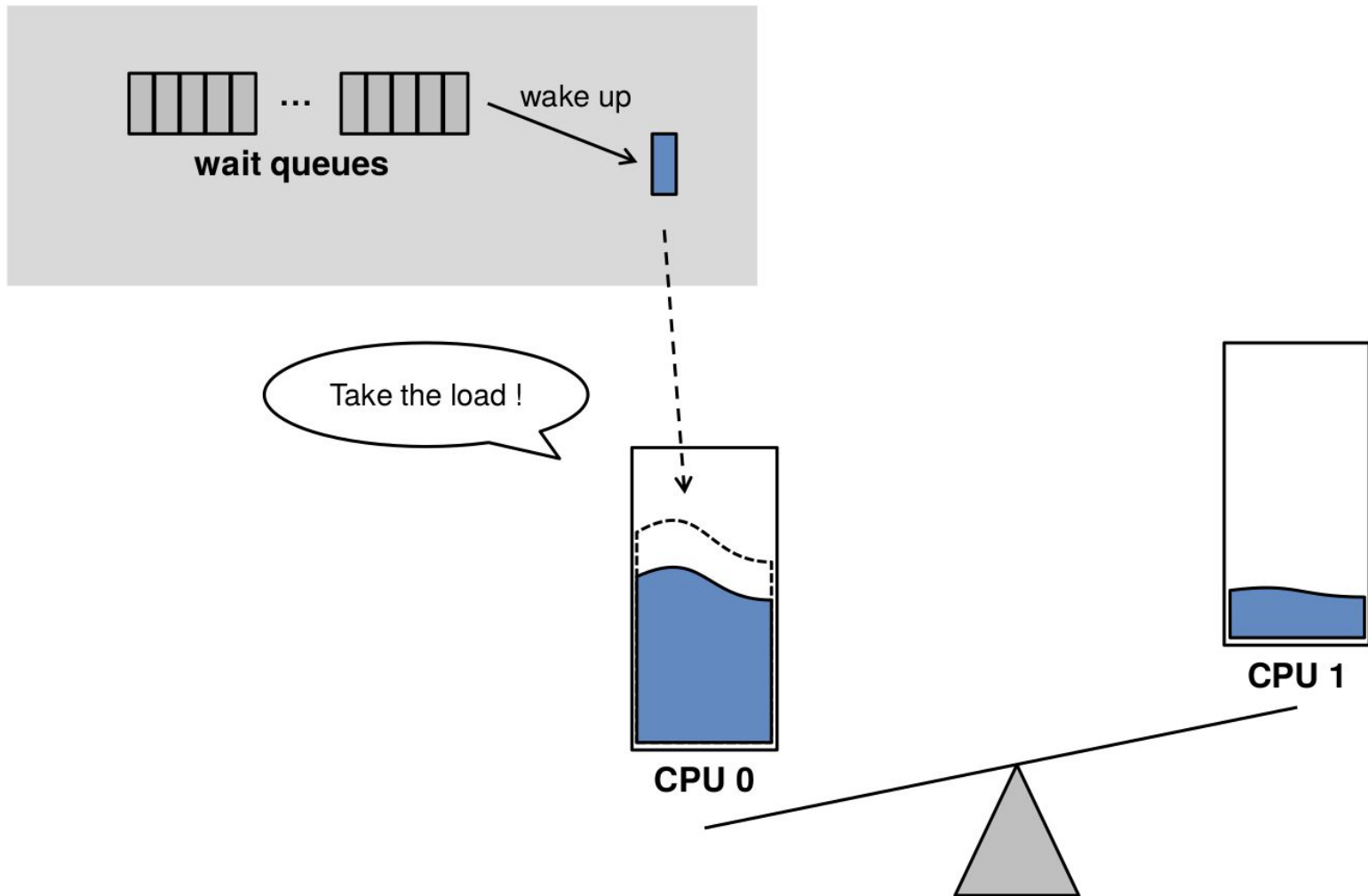




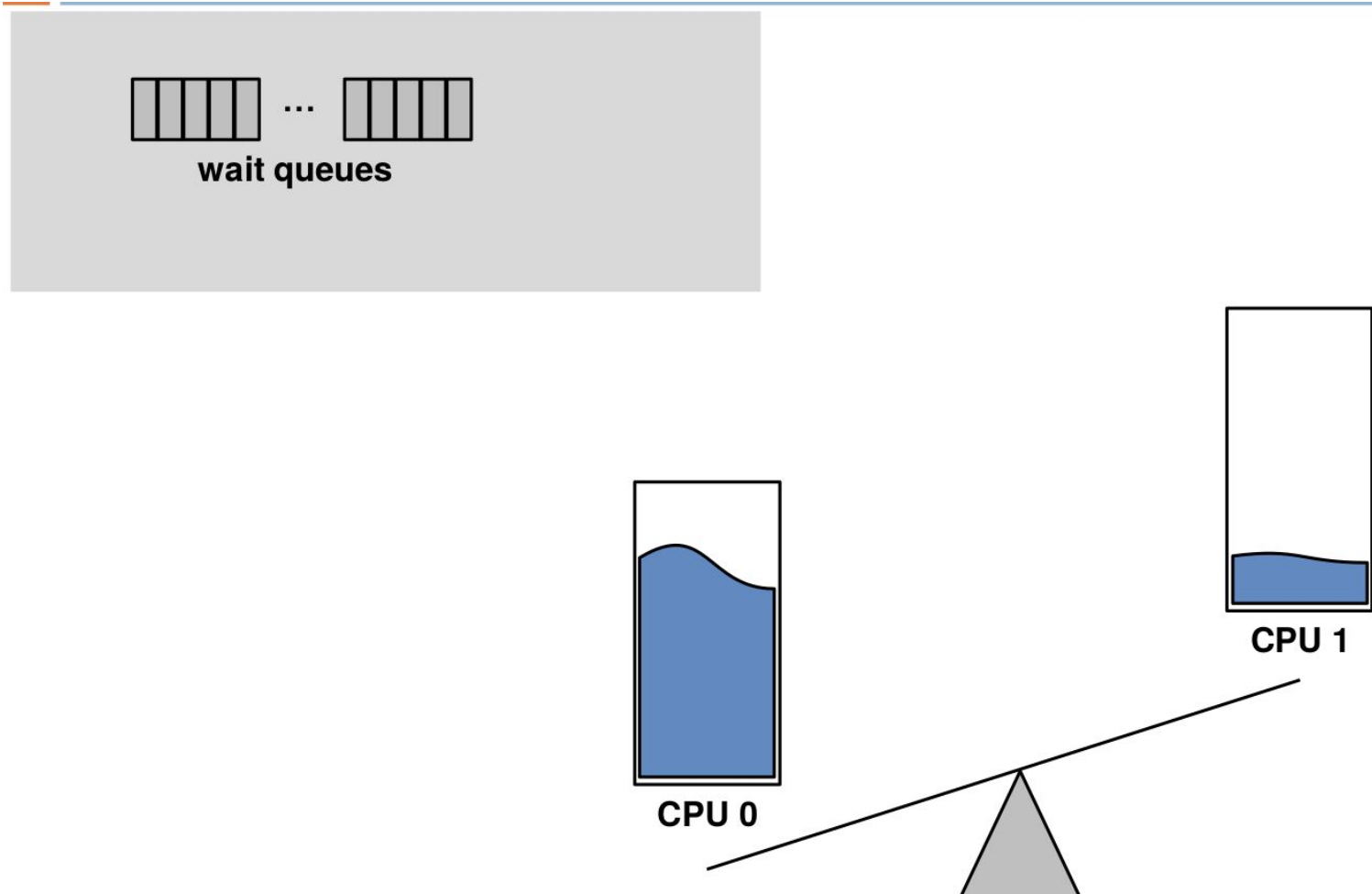
# Power Saved Load Balancing - wake-up



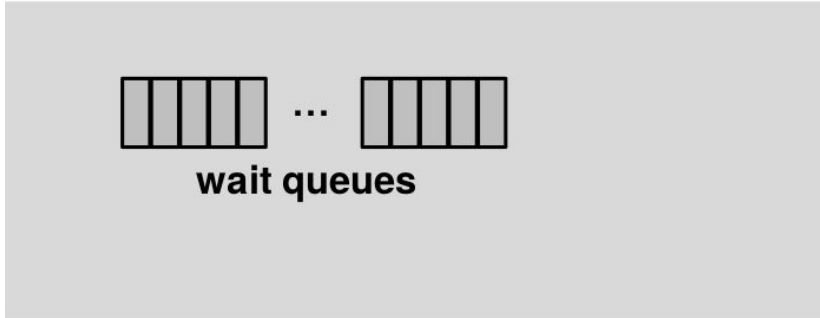
# Power Saved Load Balancing - wake-up



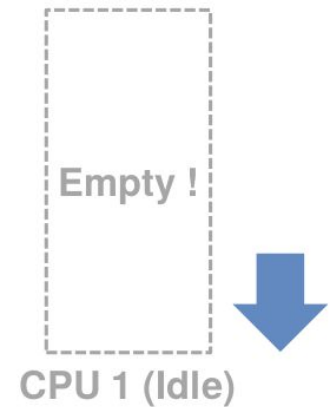
# Power Saved Load Balancing - wake-up



# Power Saved Load Balancing - wake-up

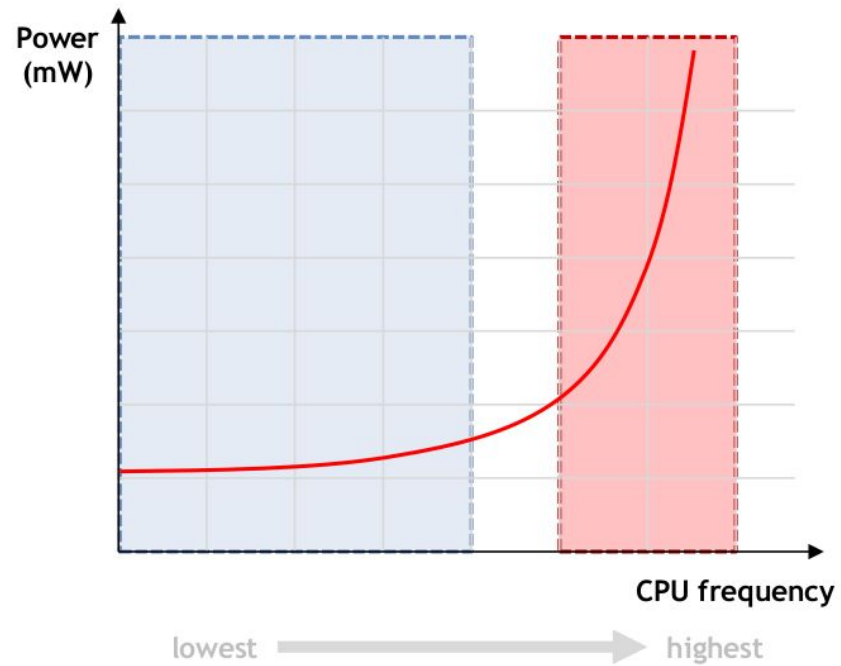


SAVE POWER !!!

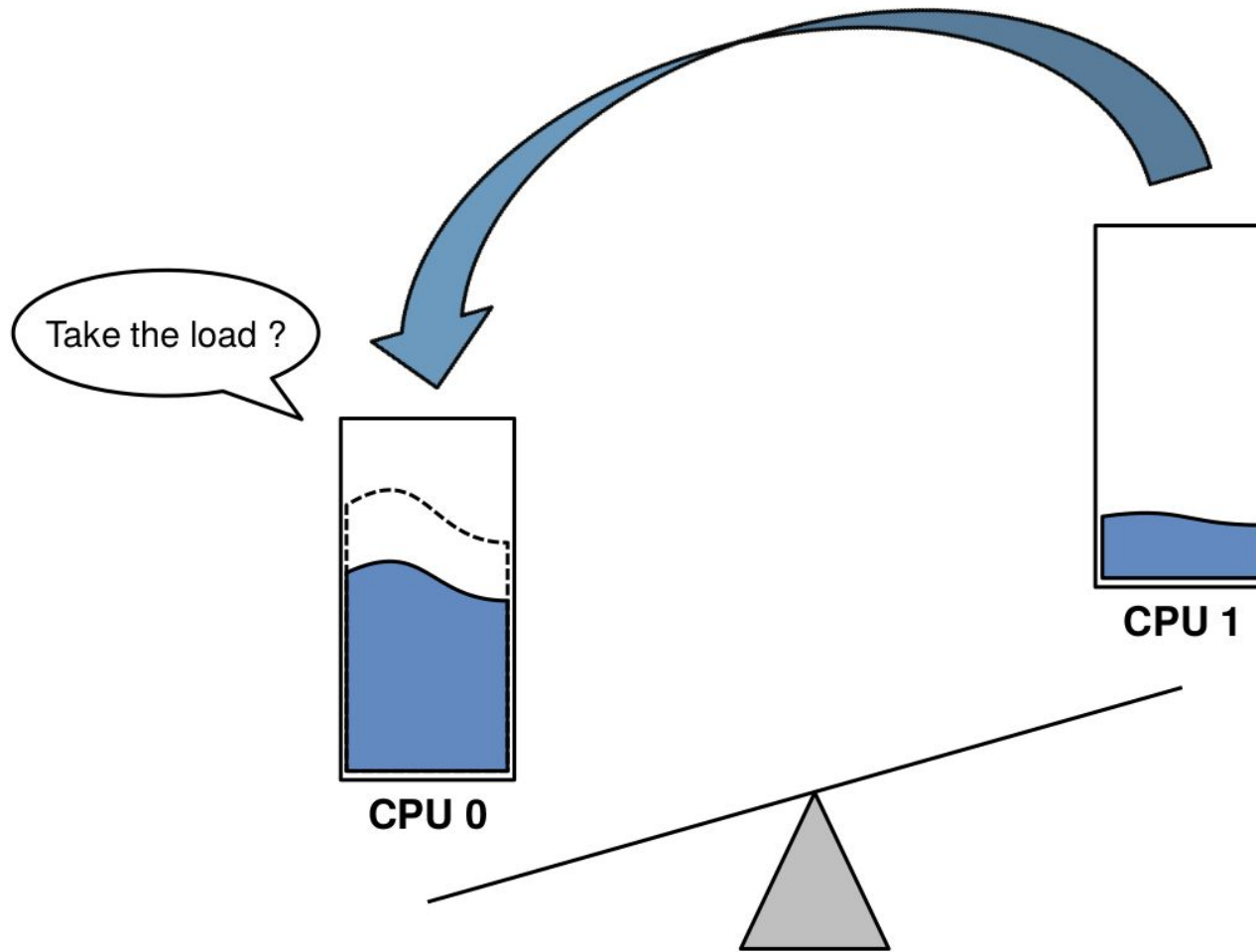


# Problem

- Frequency

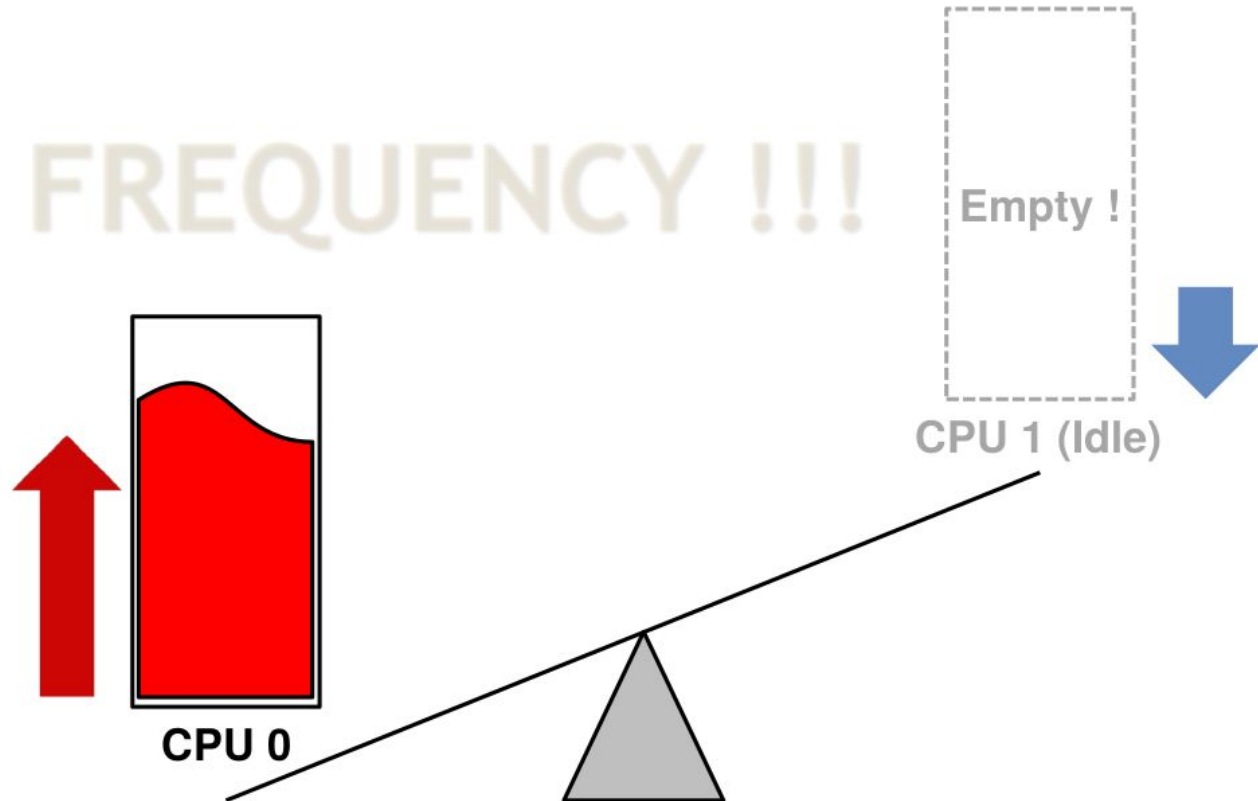


# Problem - Frequency

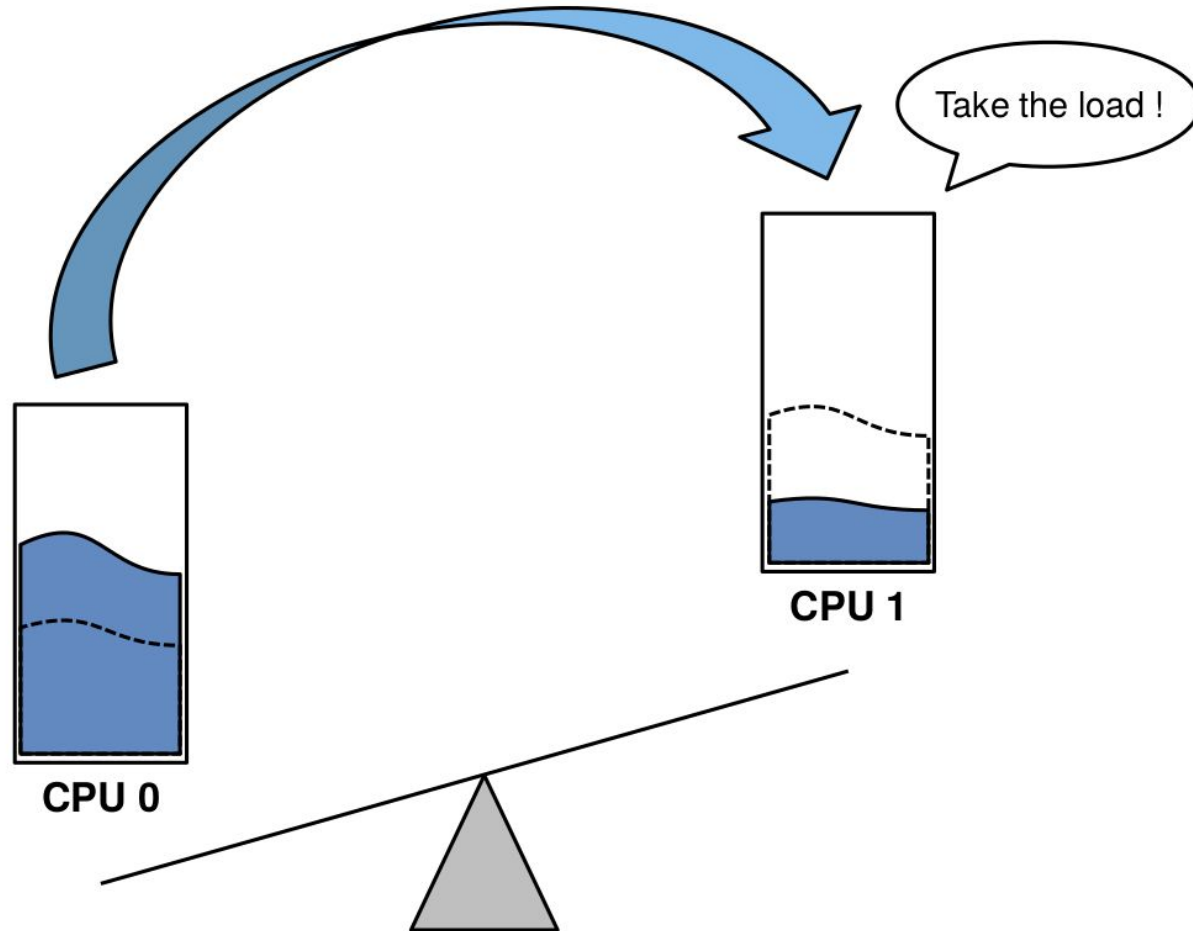


# Problem - Frequency

HIGH FREQUENCY !!!



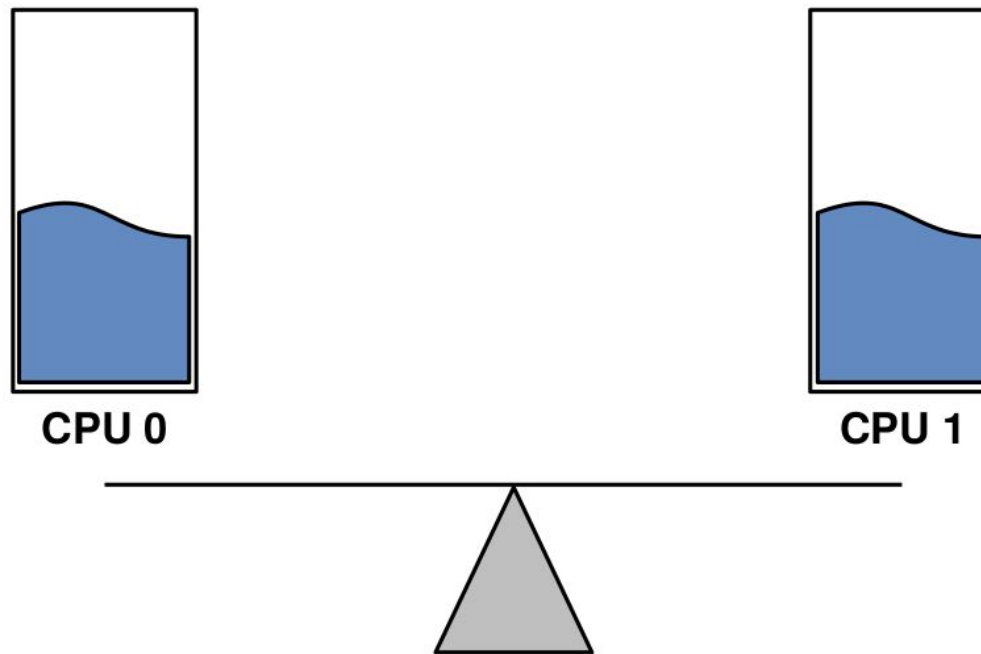
# Solution





# Solution

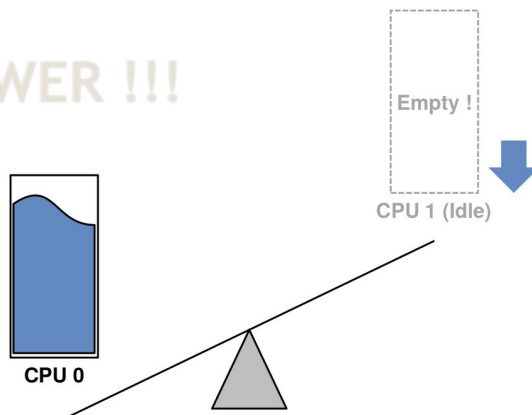
SAVE POWER !!!



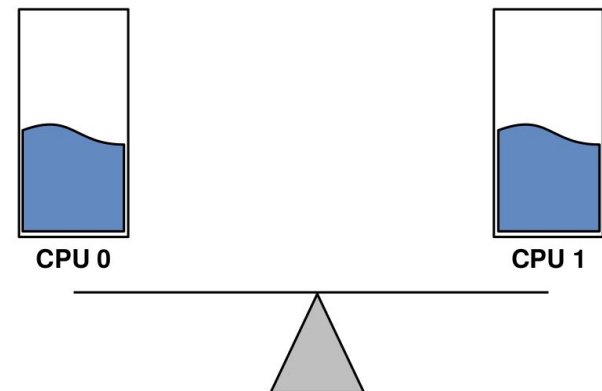
# Problem 1

- CPU idle? CPU frequency ?

SAVE POWER !!!

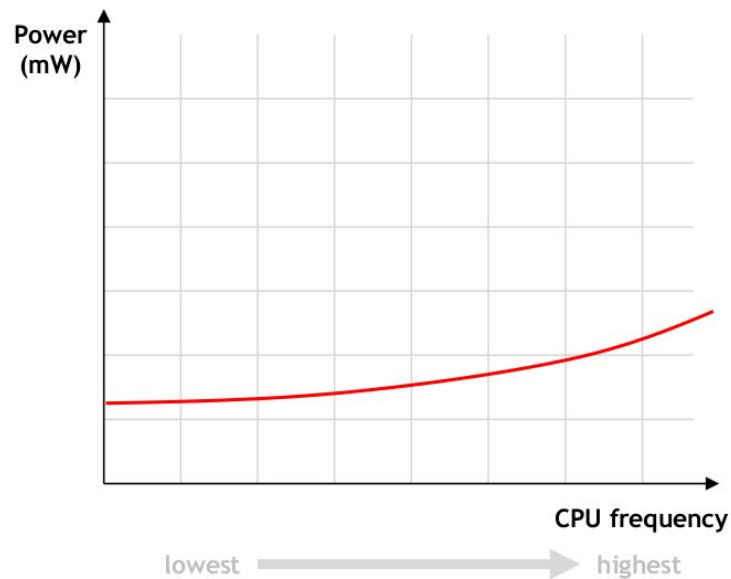


SAVE POWER !!!

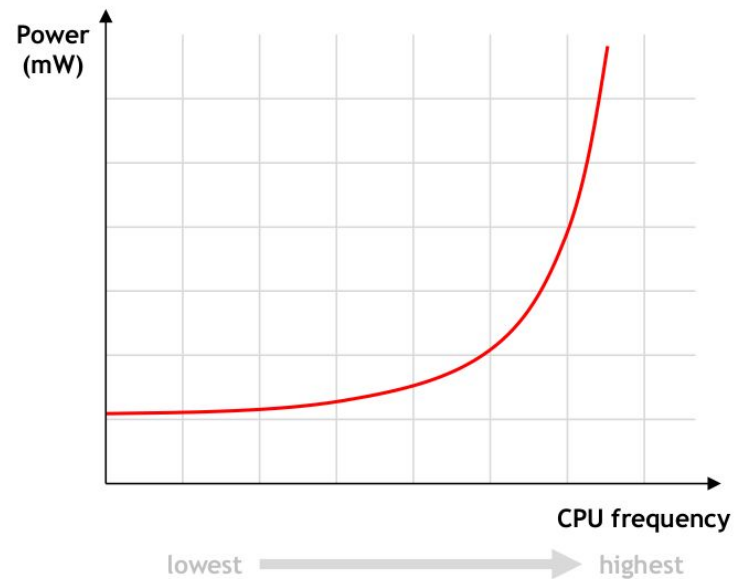


# Problem 2

- Architecture A? Architecture B?



Architecture A

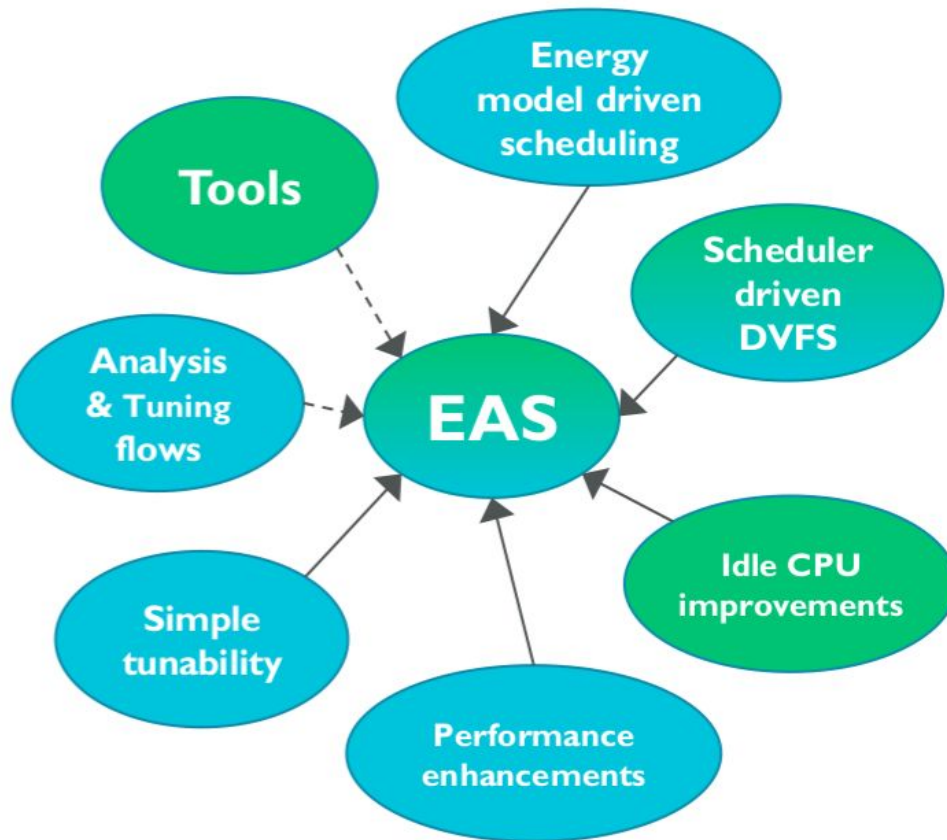


Architecture B

# Solution

- **Use Energy Cost Model**
- **For saving power consumption,**
  - Is it better to pack tasks to fewer CPUs ?
  - Is it better to spread tasks to all CPUs ?

# Solution



Linaro  
development

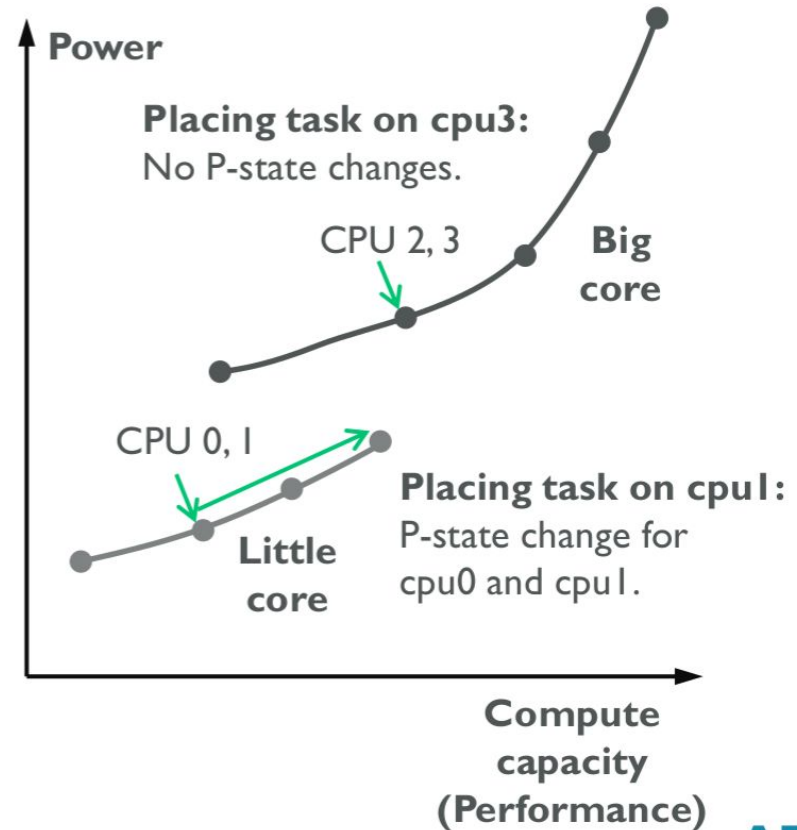
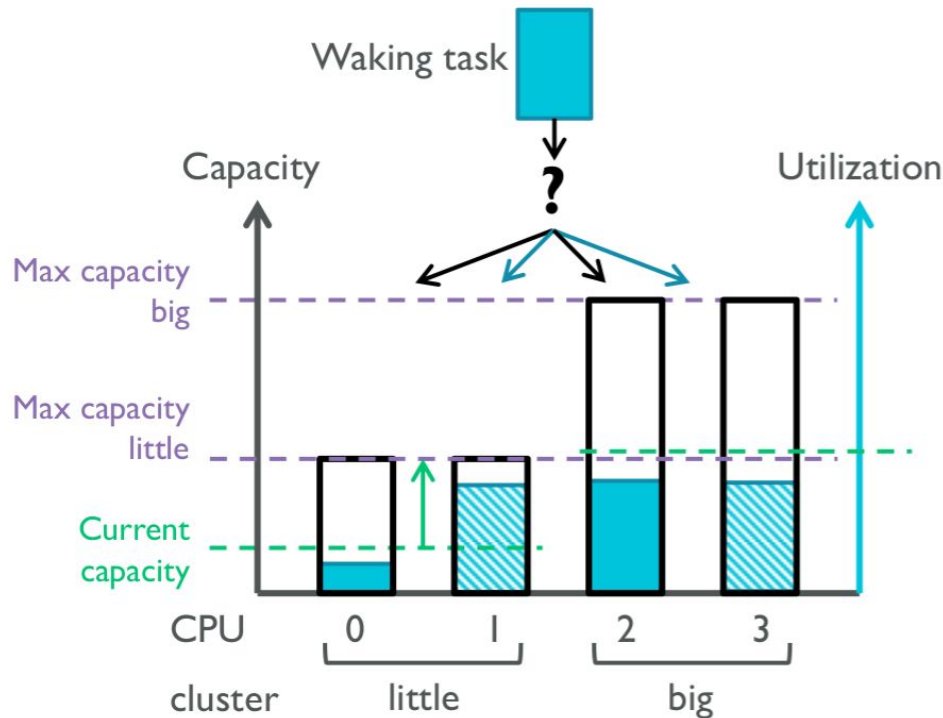
ARM  
development



# **"EAS Overview and Integration Guide"**

## **ARM**

# What is EAS – the energy model



ARM

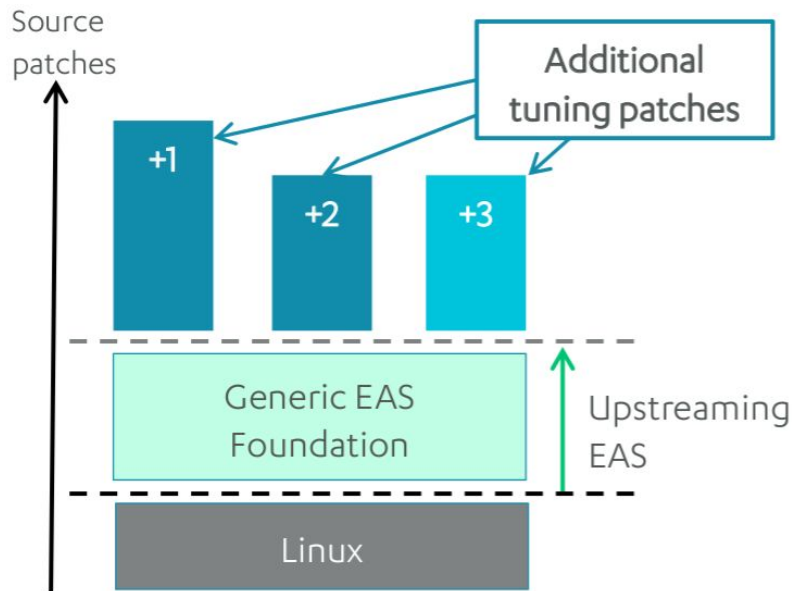


# EAS vs big.LITTLE HMP(GTS)

## EAS

New **Energy Aware Scheduling**

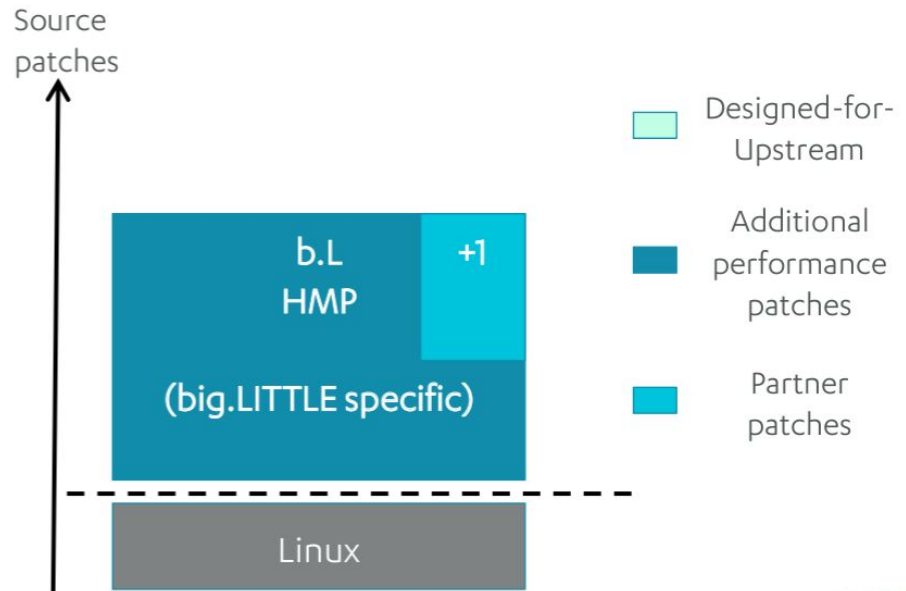
- Generic energy model based approach fits **all platforms and topologies**.
- Foundation for further enhancements.



## vs big.LITTLE HMP

Existing **Heterogeneous MP** patchset

- big.LITTLE topology only.
- Hard coded behaviors.
- In Linaro LSK kernels (not mainline).

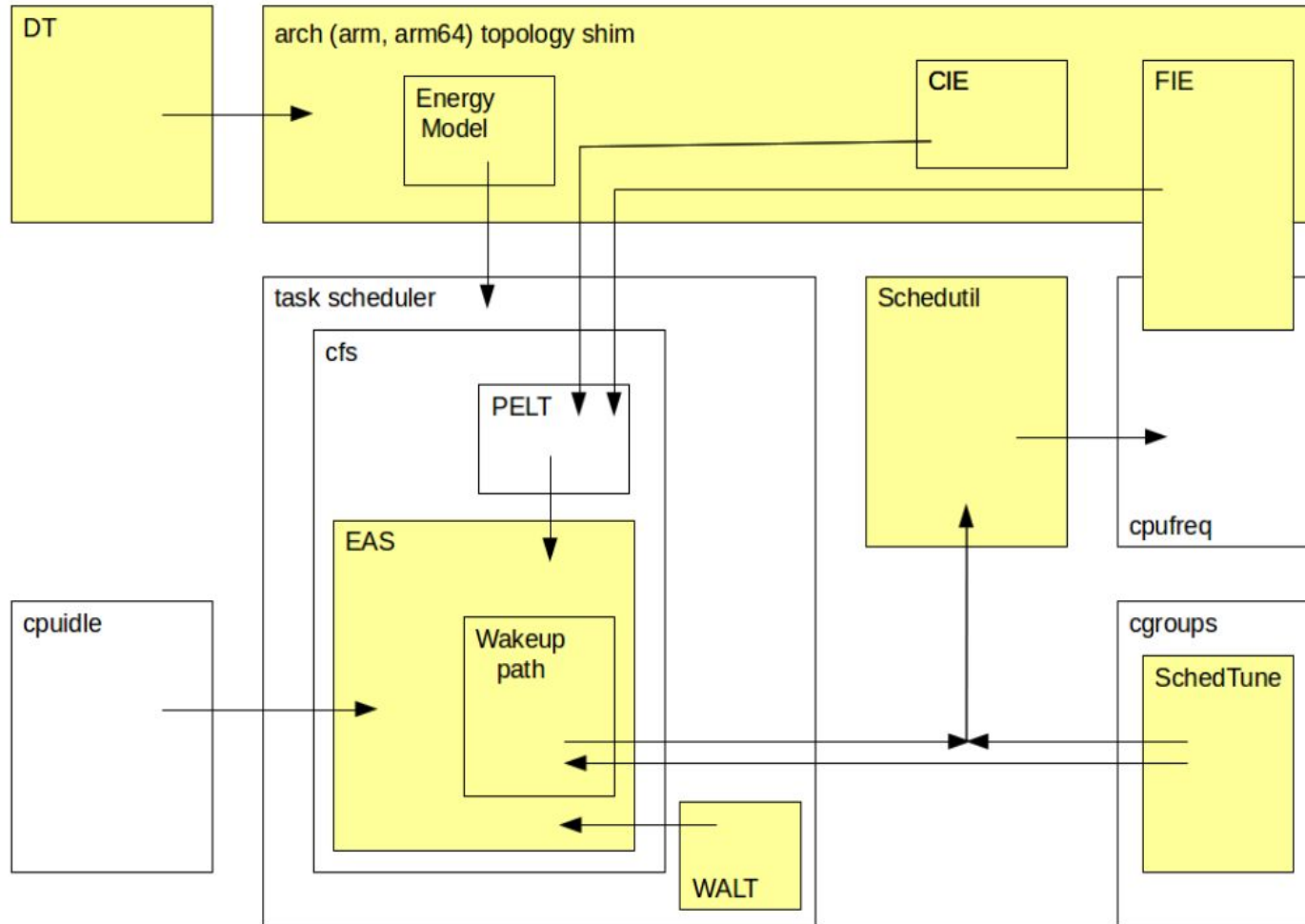


“EAS Update” 2015 september ARM

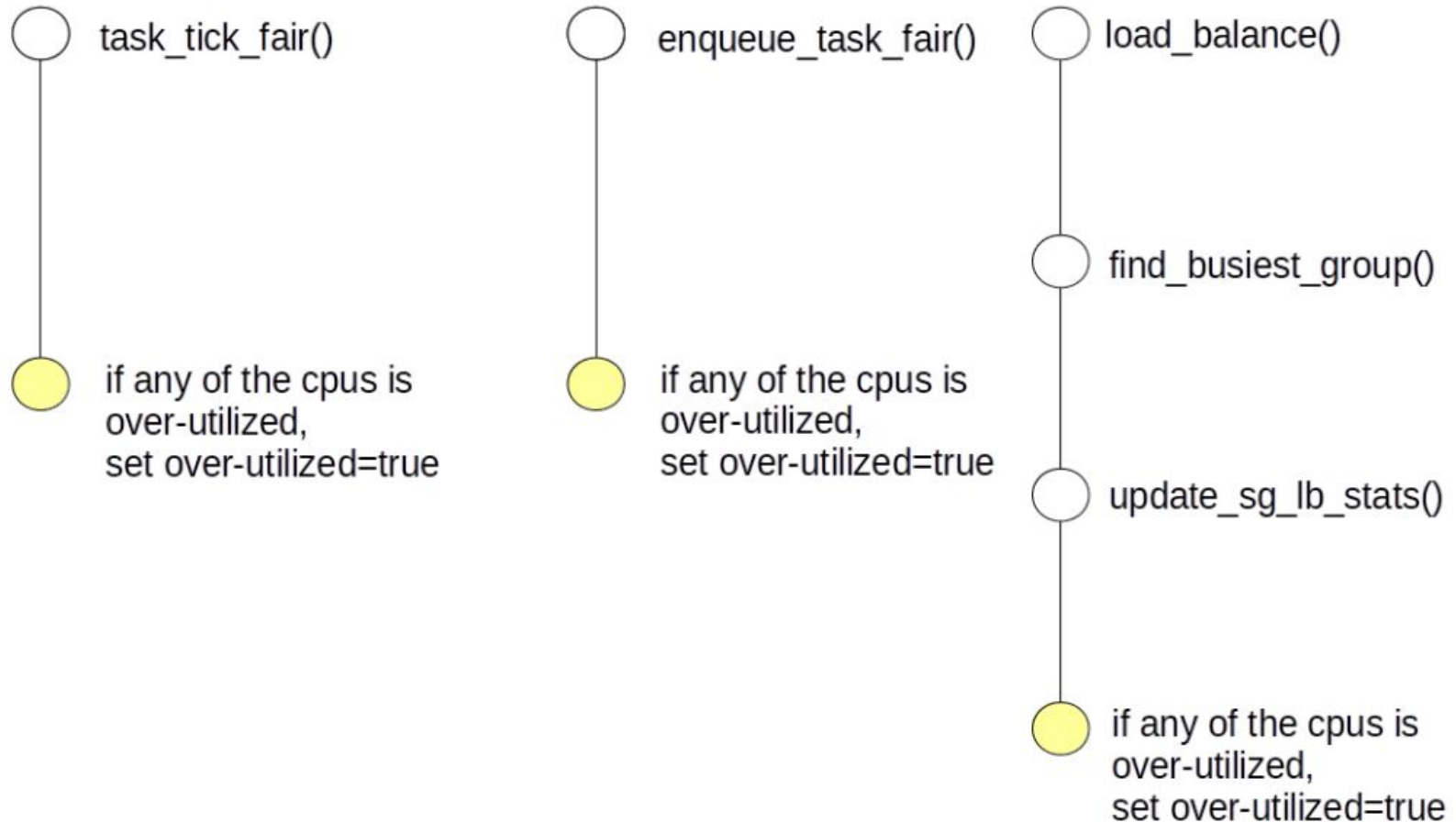
ARM



# EAS building blocks



# Over-utilized



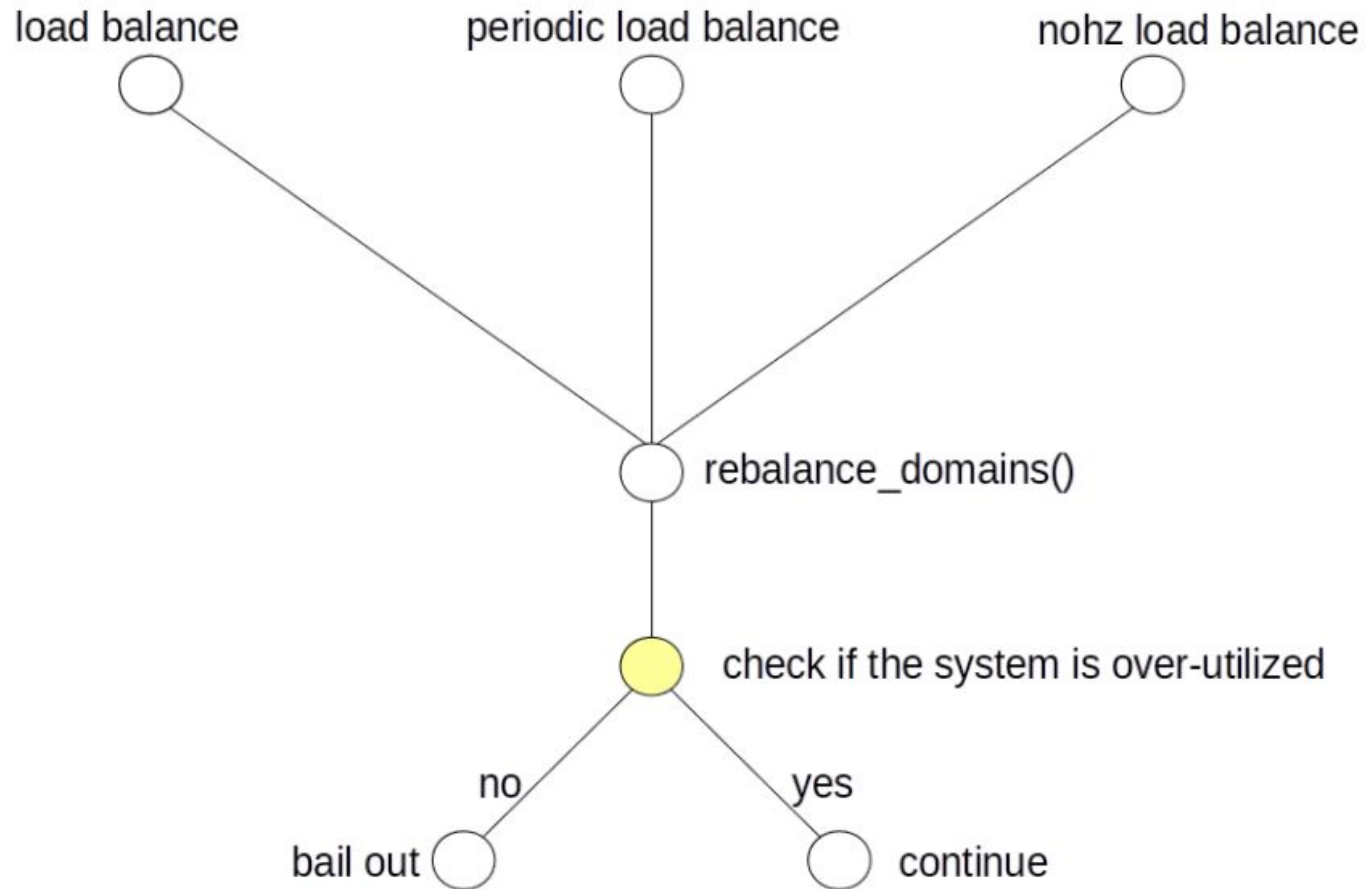
# cpu\_overutilized

```
static bool cpu_overutilized(int cpu)
{
    return (capacity_of(cpu) * 1024) < (cpu_util(cpu) * capacity_margin);
}
```

```
capacity_margin = capacity_margin = 1280; /* ~20% margin */
```

```
cpu_util = static inline unsigned long __cpu_util(int cpu, int delta)
{
    unsigned long util = cpu_rq(cpu)->cfs.avg.util_avg;
}
```

# over-utilized



# over-utilized

```
static void
enqueue_task_fair(struct rq *rq, struct task_struct
*p, int flags)
{
    [...]

    if (!se) {
        add_nr_running(rq, 1);
        if (!task_new && !rq->rd->overutilized &&
            cpu_overutilized(rq->cpu))
            rq->rd->overutilized = true;
        [...]
    }
}
```

EAS path

SMP load balance

Over tipping point

EAS path

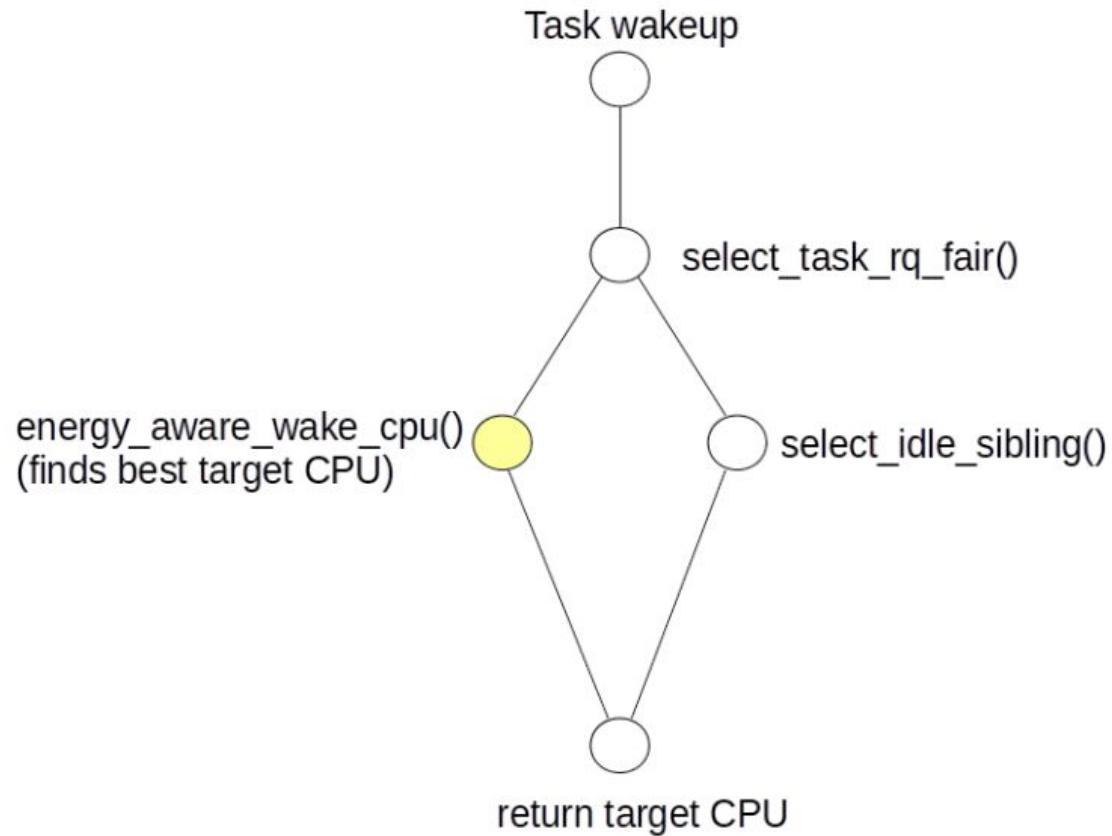
SMP load balance

```
static struct sched_group *find_busiest_group(struct lb_env
*env)
{
    if (energy_aware() && !env->dst_rq->rd->overutilized)
        goto out_balanced;
    [...]
}
```

```
static int select_task_rq_fair(struct task_struct *p, int
prev_cpu, int sd_flag, int wake_flags)
{
    [...]

    if (!sd) {
        if (energy_aware() && !cpu_rq(cpu)->rd->overutilized)
            new_cpu = energy_aware_wake_cpu(p, prev_cpu);
        else if (sd_flag & SD_BALANCE_WAKE) /* XXX always ? */
            new_cpu = select_idle_sibling(p, new_cpu);
    } else while (sd) {
        [...]
    }
}
```

# energy\_aware\_wake\_cpu



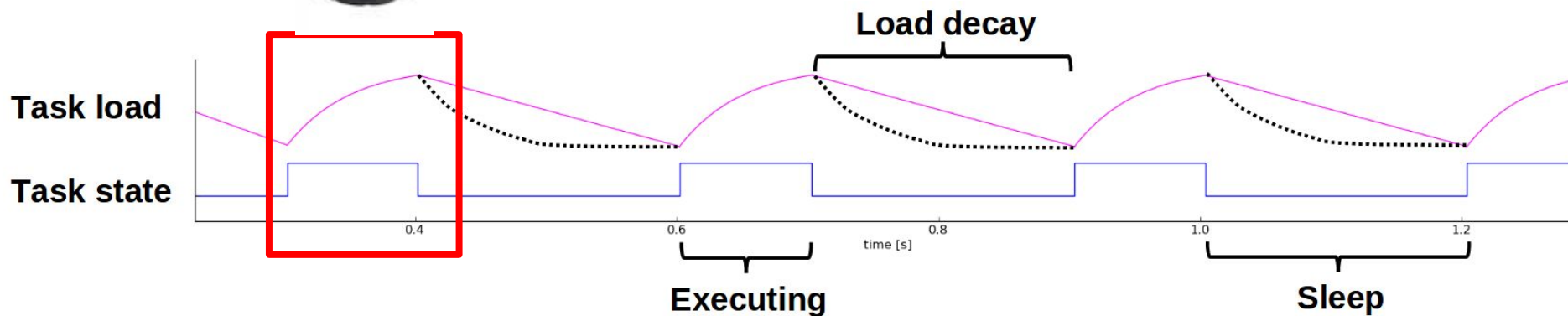


**select\_task\_rq\_fair()**

**kernel/sched/fair.c**

# Problem - PELT+EAS

- Responsiveness problem.
- Consider web browsing in mobile world.
- The load does not rise quickly enough.

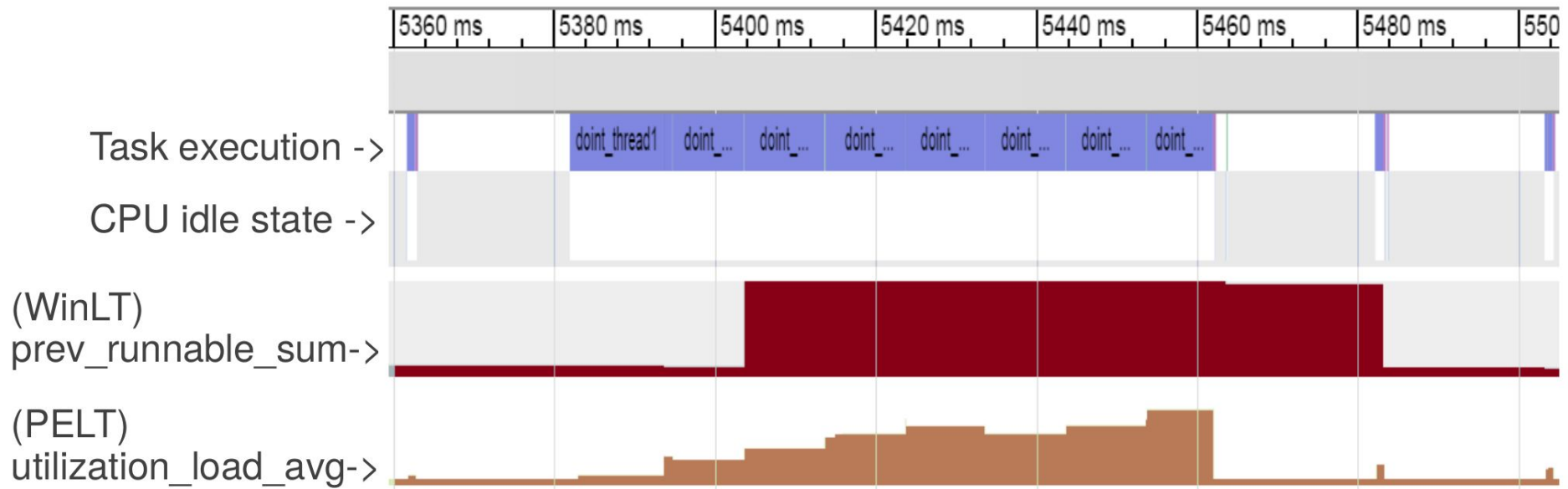




# Solution

- **Use Window Aissisted Load Tracking scheme(WALT)**

# PELT and WinLT

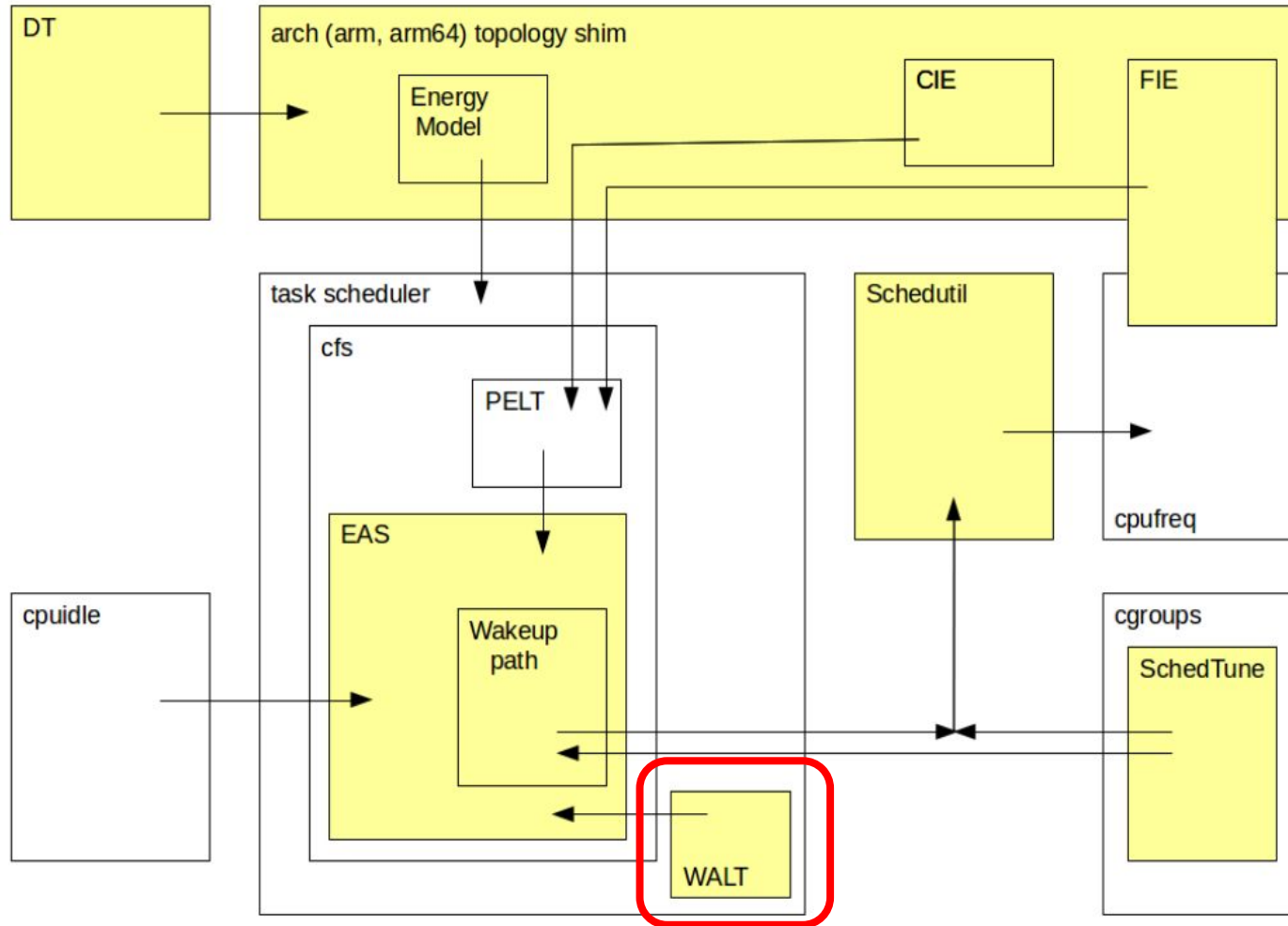


<https://www.slideshare.net/linaroorg/bkk16208-eas>

# Solution

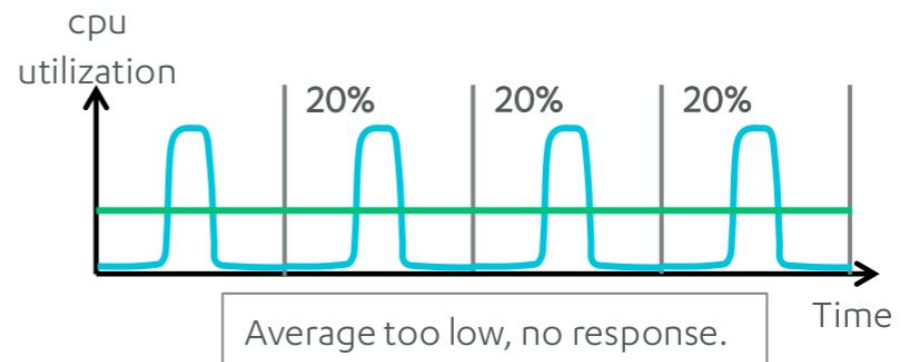
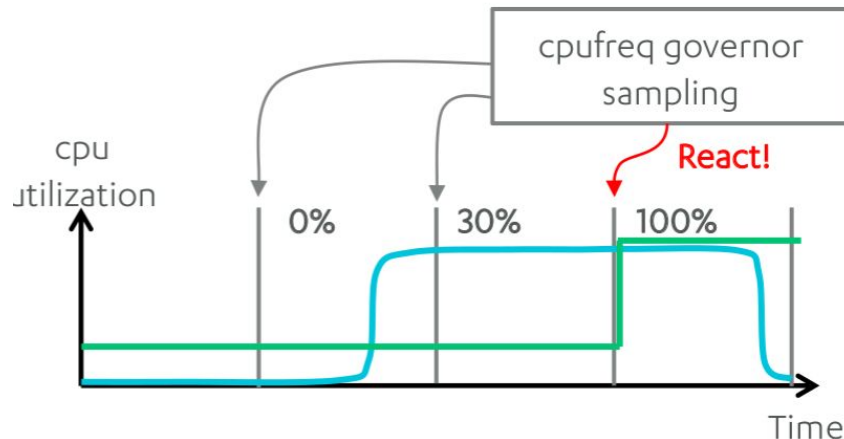
- **Use Window Aissisted Load Tracking scheme(WALT)**
- **WALT**
  - Problem : Lock + Cache coherency
- **Solution**
  - new PELT Load tracking algorithm

# EAS building blocks



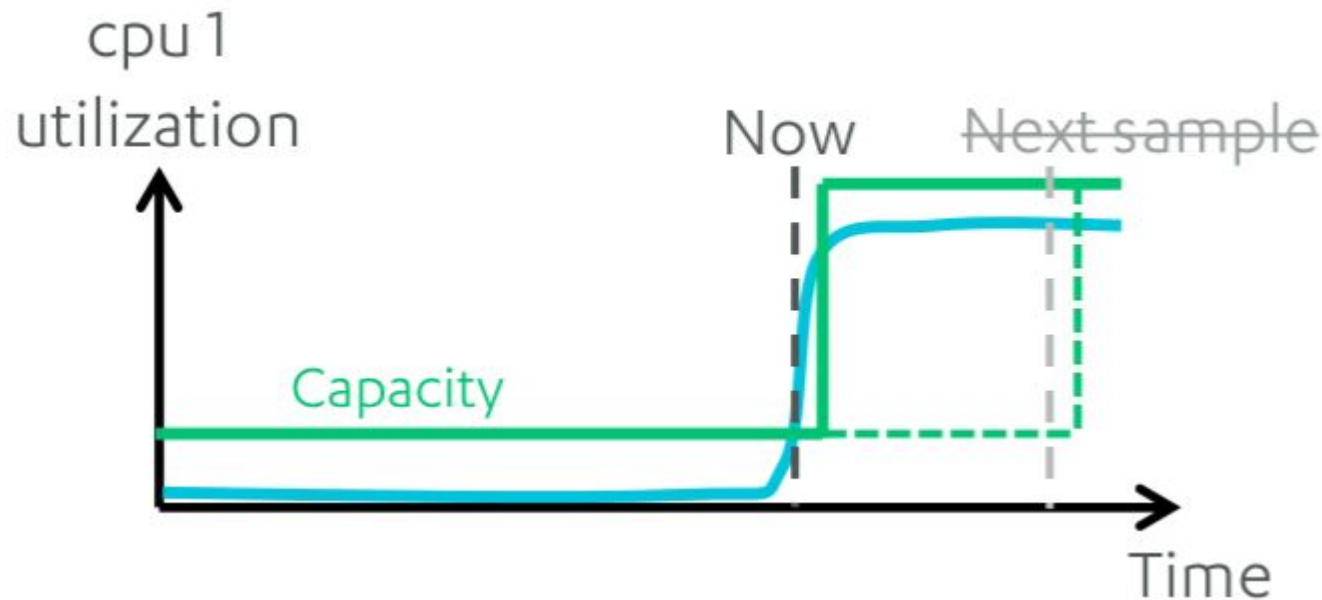
# Problem - Timer-based DVFS

- **Problem**
  - Sampling based governors are slow to respond and hard to tune.
- **Sampling too fast**
  - Freq changes for small utilization spikes.
- **Sampling too slow**
  - Average too low

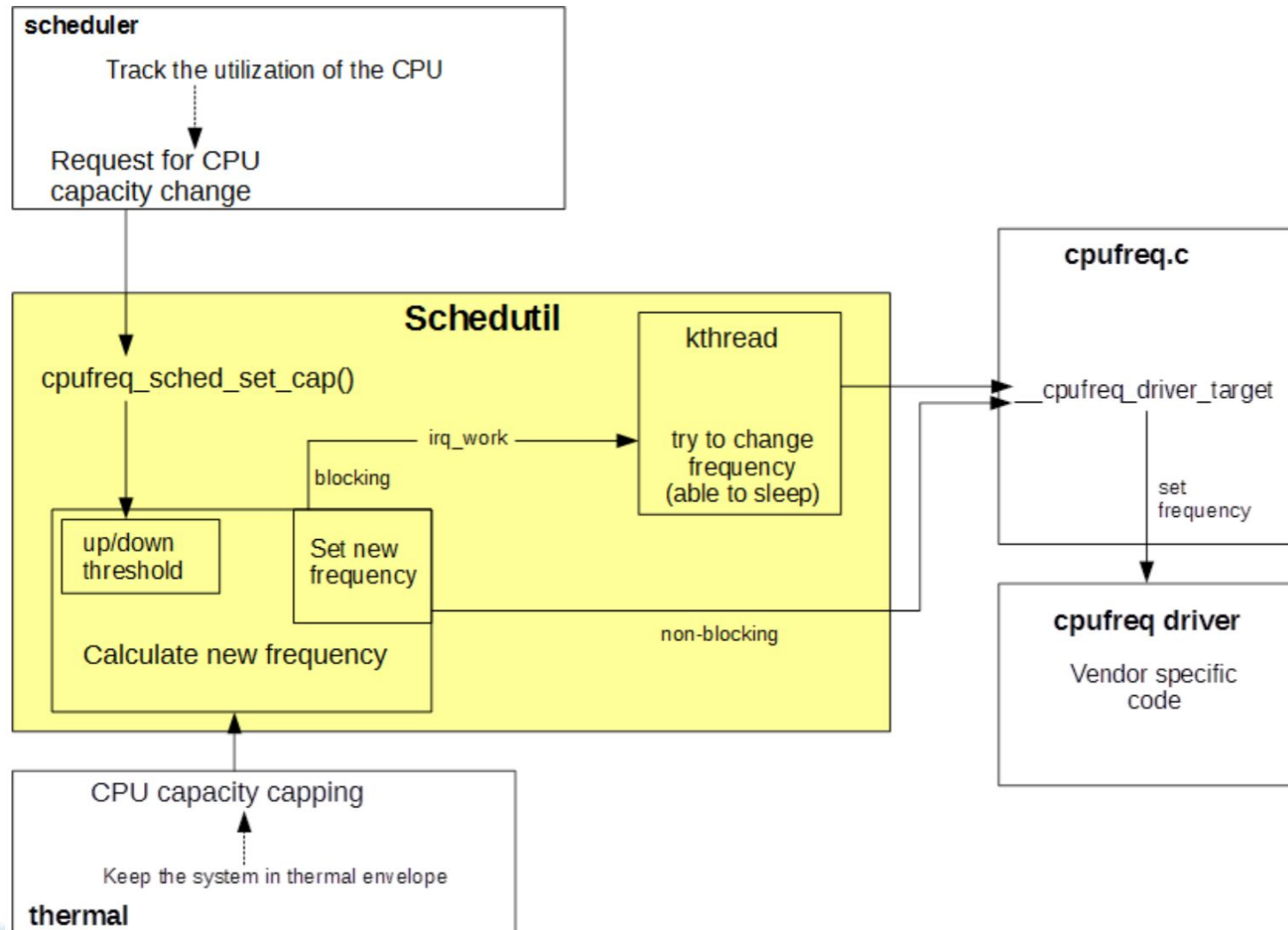


# Solution

- Use scheduler-based DVFS
- Scheduler task utilization tracking DVFS
- Immediately when CPU utilization changes
  - improved responsiveness.



# Schedutil





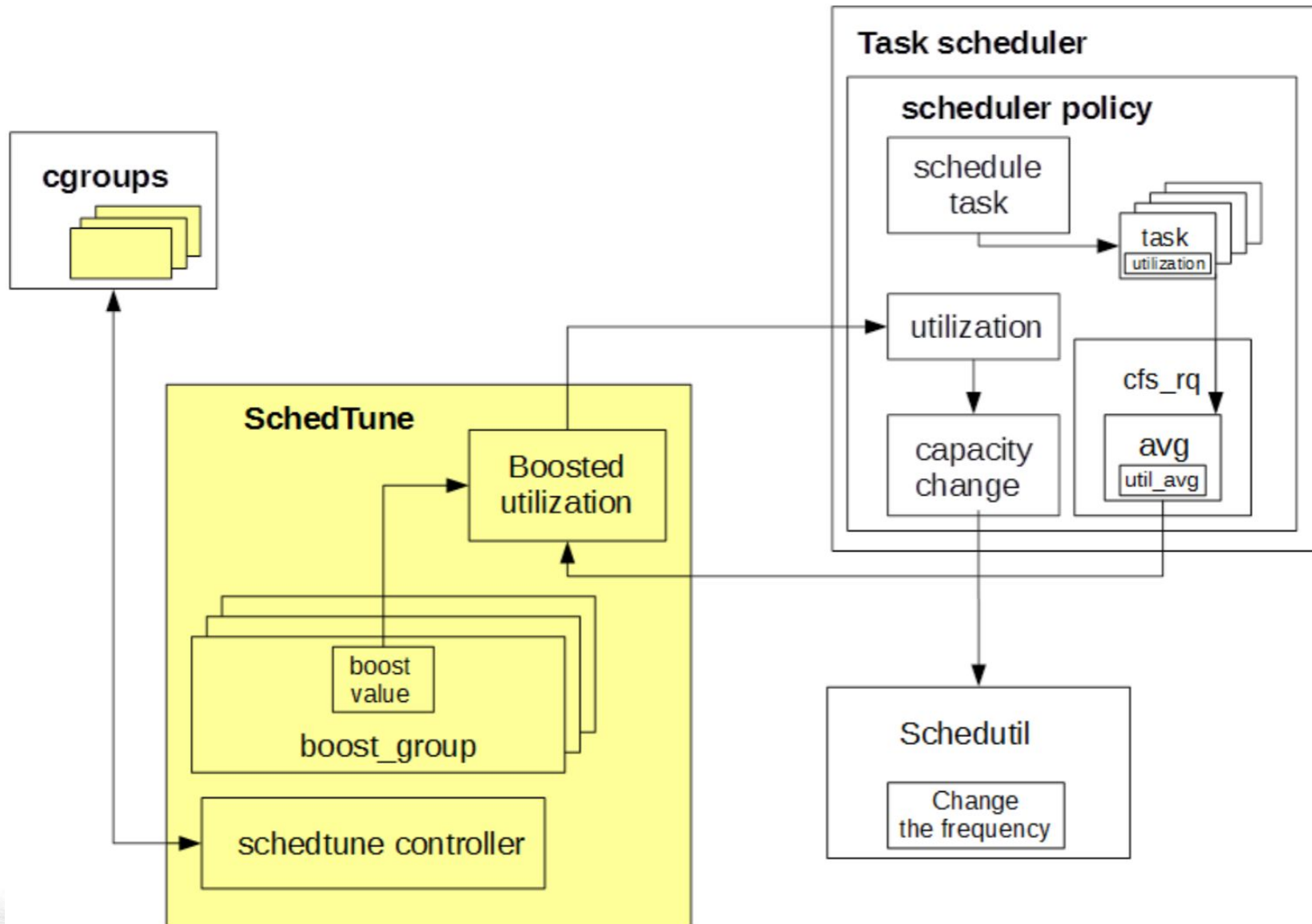
# Problem - Android

- **Sometimes response time is crucial to the users.**
- **Scheduler**
  - Cannot give user's experience information.

# Solution

- **Use SchedTune**
- **A way of to provide user's experience information.**
- **Implemented as a control-group controller.**
- **Each control group has a single tunable knob.**
  - `schedtune.boost`

# SchedTune



# Create energy-aware scheduler tuning nodes

```
mkdir /dev/stune
mount cgroup none /dev/stune schedtune
mkdir /dev/stune/foreground
mkdir /dev/stune/background
mkdir /dev/stune/top-app
chown system system /dev/stune
chown system system /dev/stune/foreground
chown system system /dev/stune/background
chown system system /dev/stune/top-app
chown system system /dev/stune/tasks
chown system system /dev/stune/foreground/tasks
chown system system /dev/stune/background/tasks
chown system system /dev/stune/top-app/tasks
chmod 0664 /dev/stune/tasks
chmod 0664 /dev/stune/foreground/tasks
chmod 0664 /dev/stune/background/tasks
chmod 0664 /dev/stune/top-app/tasks
```

# Reference

- <https://pdos.csail.mit.edu/6.828/2016/schedule.html>
- <http://web.mit.edu/6.033>
- <http://www.rdrop.com/~paulmck/>
- "Is Parallel Programming Hard, And If So, What Can You Do About It?"
- Davidlohr Bueso. 2014. Scalability techniques for practical synchronization primitives. *Commun. ACM* 58  
<http://queue.acm.org/detail.cfm?id=2698990>
- "CPUFreq and The Scheduler Revolution in CPU Power Management", Rafael J. Wysocki
- <https://sites.google.com/site/embedwiki/oses/linux/pm/pm-qos>
- <https://intl.aliyun.com/forum/read-916>
- User-level threads : co-routines  
<http://www.gamedevforever.com/291>  
[https://www.youtube.com/watch?v=YYtzQ355\\_Co](https://www.youtube.com/watch?v=YYtzQ355_Co)
- Scheduler Activations
  - <https://cgi.cse.unsw.edu.au/~cs3231/12s1/lectures/SchedulerActivations.pdf>
- [https://en.wikipedia.org/wiki/FIFO\\_\(computing\\_and\\_electronics\)](https://en.wikipedia.org/wiki/FIFO_(computing_and_electronics))
- <http://jake.dothome.co.kr/>
- <http://www.linuxjournal.com/magazine/completely-fair-scheduler?page=0.0>
- [https://www2.cs.uic.edu/~jbell/CourseNotes/OperatingSystems/6\\_CPU\\_Scheduling.html](https://www2.cs.uic.edu/~jbell/CourseNotes/OperatingSystems/6_CPU_Scheduling.html)
- "Energy Aware Scheduling", Byungchul Park, LG Electronic
- "Update on big.LITTLE scheduling experiments", ARM
- "EAS Update" 2015 september ARM
- "EAS Overview and Integration Guide", ARM TR
- "Drowsy Power Management", Matthew Lentz, SOSP 2015
- <https://www.slideshare.net/nanik/learning-aosp-android-hardware-abstraction-layer-hal>
- <https://www.youtube.com/watch?v=oTGQXqD3CNI>
- <https://www.youtube.com/watch?v=P80NcKUKpuo>
- <https://lwn.net/Articles/398470/>
- "SCHED\_DEADLINE: It's Alive!", ARM, 2017