

Literature Review *of* HMP Scheduling

□ big.LITTLE technology

“big” cores should provide maximal computing power

“LITTLE” cores are designed for maximum power efficiency

Peter Hu

Huawei Technologies Research & Development

CPU Team

19 June 2023



Table of Contents with *references*

Existing Linux CFS and EAS implementation

Ref: EAS Overview and Integration Guide, arm, 2018

<https://docs.kernel.org/scheduler>

Qualcomm HMP scheduler

Ref: EAS – Energy Aware Scheduler An unbiased look, Vitaly Wool, Konsulko Group

Multimedia, IEEE & Samsung Electronics

saves system-wide (not just CPU) energy consumption by 8.9 percent

Ref: Enhancing Energy Efficiency of Multimedia Applications in Heterogeneous Mobile Multi-Core Processors (2017)

Abstraction Allocates multimedia applications to the small cores and non-multimedia applications to the big cores.

WAEAS, TSINGHUA SCIENCE AND TECHNOLOGY

improves 5-8% when the workload was high

Ref: An Optimization Scheme of EAS Scheduler for Wearable Applications (2021)

BES-WALT Basic exponential smoothing-based WALT algorithm, Energy-aware CPU selection algorithm,

Abstraction Adjusts Sched Tune

Latency-Sensitive Task: decrease searching speed to get minimum performance capacity and the lowest idle state

Non-Latency-Sensitive Task: not choose a lower utilization, but “task packing strategy” (on the little cores)

Batch processing strategy,

Abstraction Adjusts *select idle sibling()*

Record and centralize on fewer cores

Cluster-based load balancing (overutilized)

Abstraction Reduces the meaningless task migration by local load balancing
(when having many aperiodic high-load applications)

Table of Contents with *references*

Learning EAS, LG Electronics

improves power consumption by 2.8% - 7.8%

Ref: Performance Improvement of Linux CPU Scheduler Using Policy Gradient Reinforcement Learning for Android Smartphones (2020)

Abstraction: Adjusts the *TARGET_LOAD* used to set the CPU frequency and the *sched migration cost* used as the task migration criteria
Using Policy reinforcement learning dealing with **workload** or the **ratio of sleep and running states** changes.

Other *references*

Wider coverage: [Energy-Aware Scheduling for High-Performance Computing Systems: A Survey \(2023\)](#)

Terminology

- Completely Fair Scheduler models an "ideal, precise multi-tasking CPU" on real hardware.

- CFS maintains the amount of time provided to a given task to determine if it needs balancing
 - the smaller amount of time a task has been permitted access to the processor — the higher its need for the processor is
- CFS maintains a time-ordered red-black tree, sorting tasks in ascending order by CPU bandwidth received
 - Instead of run queues as did predecessors
 - Guarantees $O(\log(N))$
- The leftmost task off the red-black tree is picked up next
 - It has the least spent execution time
 - So that task gets the CPU to restore balance (fairness)
- Considers all CPUs to be the same
 - Works very well in SMP systems
 - Does not work in more complicated cases

Source: EAS – Energy Aware Scheduler *An unbiased look*, Vitaly Wool, Konsulko Group

Original CFS policy operates at **full** system utilization, while EAS operates when the system has low/medium total utilization, as EAS offers no energy benefit when all CPUs are overutilized

<https://docs.kernel.org/scheduler/sched-design-CFS.html>

<https://docs.kernel.org/scheduler/sched-energy.html>

Existing Linux CFS and EAS implementation

Aim: Energy-Efficiency and Performance.
Fairness

Terminology

- Energy = [joule] (resource like a battery on powered devices)
- Power = energy/time = [joule/second] = [watt]
- CPU Capacity = $\text{work_per_hz(cpu)} * \text{max_freq(cpu)}$

Millions of Instructions Per Second (MIPS)

*the **performance** a CPU can reach, normalized against the most performant CPU in the system.*

Heterogeneous systems have asymmetric CPU capacity.

The goal of EAS is to minimize energy, while still getting the job done.

maximize: *performance [inst/s] / Power [W]*

minimize: *Energy [J] / instruction*

EAS:

To break the tie between several good CPU candidates and pick the one that is predicted to yield the **best energy consumption** without harming the system's throughput when deciding which a task should run during wake-up.

Rely on specific elements about the platform's topology, the 'capacity' of CPUs, and their respective energy costs.

<https://docs.kernel.org/scheduler/sched-design-CFS.html>

<https://docs.kernel.org/scheduler/sched-energy.html>

Device Tree Source File / Energy Model

```
energy-costs {
    CPU_COST_0: core-cost0 {
        busy-cost-data = <
            417  168
            579  251
            744  359
            883  479
            1024 616
        >;
        idle-cost-data = <
            15
            0
        >;
    };
    CPU_COST_1: core-cost1 {
        busy-cost-data = <
            235 33
            302 46
            368 61
            406 76
            447 93
        >;
        idle-cost-data = <
            6
            0
        >;
    };
};
```

```
CLUSTER_COST_0: cluster-cost0 {
    busy-cost-data = <
        417  24
        579  32
        744  43
        883  49
        1024 64
    >;
    idle-cost-data = <
        65
        24
    >;
};
CLUSTER_COST_1: cluster-cost1 {
    busy-cost-data = <
        235 26
        303 30
        368 39
        406 47
        447 57
    >;
    idle-cost-data = <
        56
        17
    >;
};
```

Existing Linux CFS and EAS implementation

Terminology

Load Tracking Mechanism

Per- Entity Load Tracking (PELT) - focused only on a single class - CFS

Let L_i designate the entity's load contribution in period p_i

Then the total load is
$$L = L_0 + L_1q + L_2q^2 + L_3q^3 + \dots$$

q is the decay factor

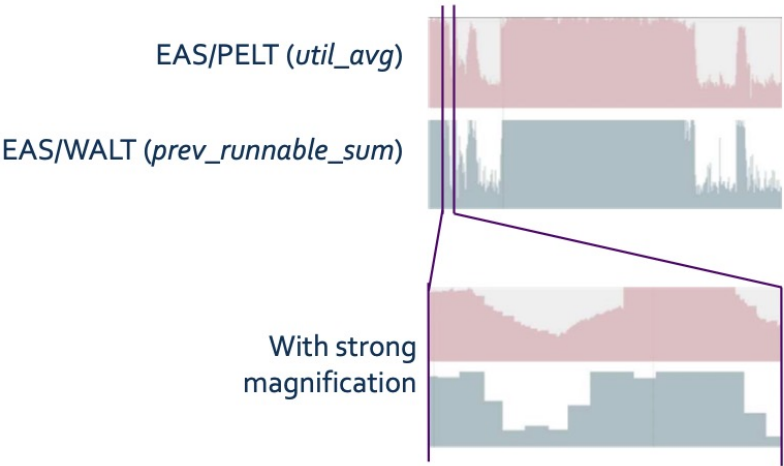
	PELT	WALT
Load tracking	Load is accounted using a geometric series	Load is accounted with a policy that observes past N windows
Blocked load/utilization tracking	Load is decayed as part of a runqueue statistic when the task is blocked	Blocked load contribution is removed from runqueue sum/average statistics.
Blocked load restoration	Runqueue statistics include blocked load/utilization at all times	Load contribution is restored to RQ statistics when the task becomes runnable again.

Window Assist Load Tracking (WALT) - EAS

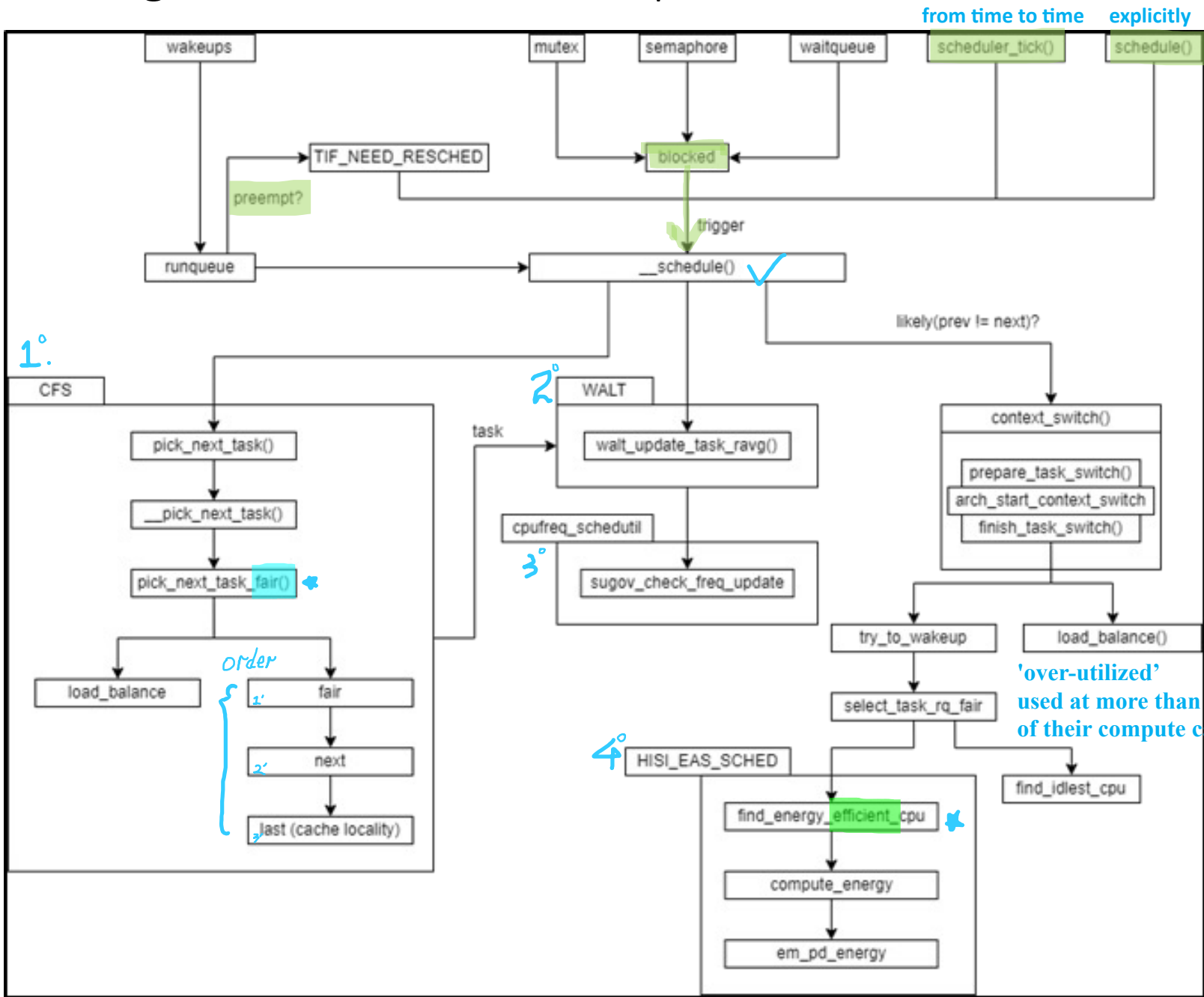
faster reaction times when the behaviour of tasks changes.
uses periodic calculations that are synchronized across all of the run queues, attempting to track the behaviour of **all** scheduling classes.

- ✓ the decisions can be made based on the information about the full state of the running system.
- ✗ additional locking complexity and some additional delays in other pathways

<https://docs.kernel.org/scheduler/sched-design-CFS.html>
<https://docs.kernel.org/scheduler/sched-energy.html>



Existing Linux CFS and EAS implementation



CPU capacity
a measure of the **performance** a CPU can reach, normalized against the most performant CPU in the system.

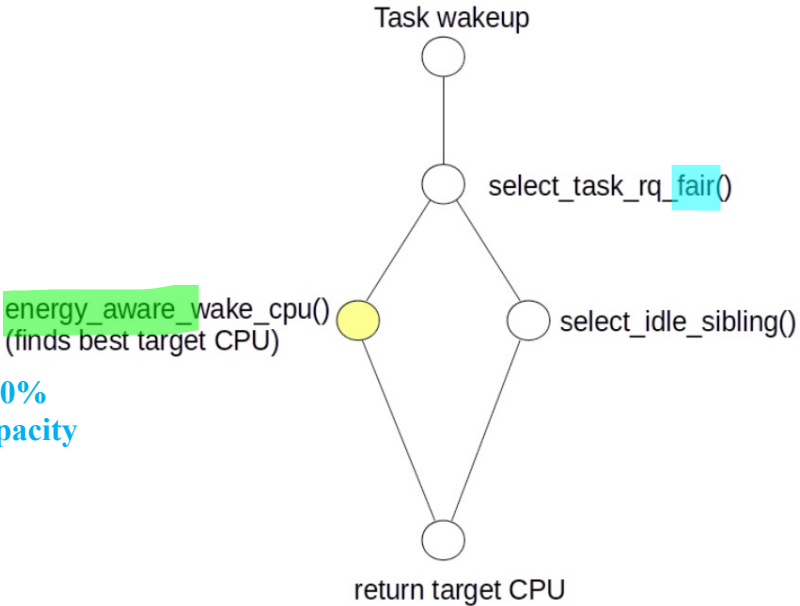
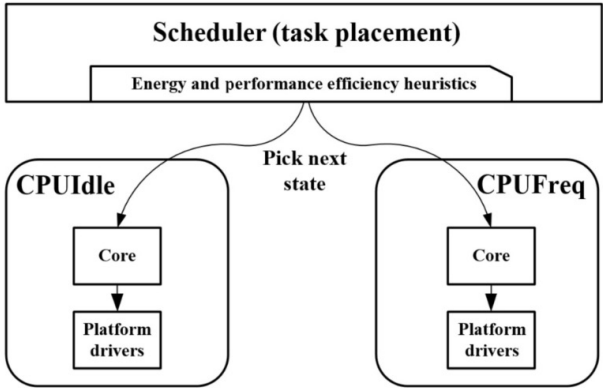


Figure 2 Task wakeup path modifications. (The yellow dots represent EAS functionality)

1.8 EAS Overview

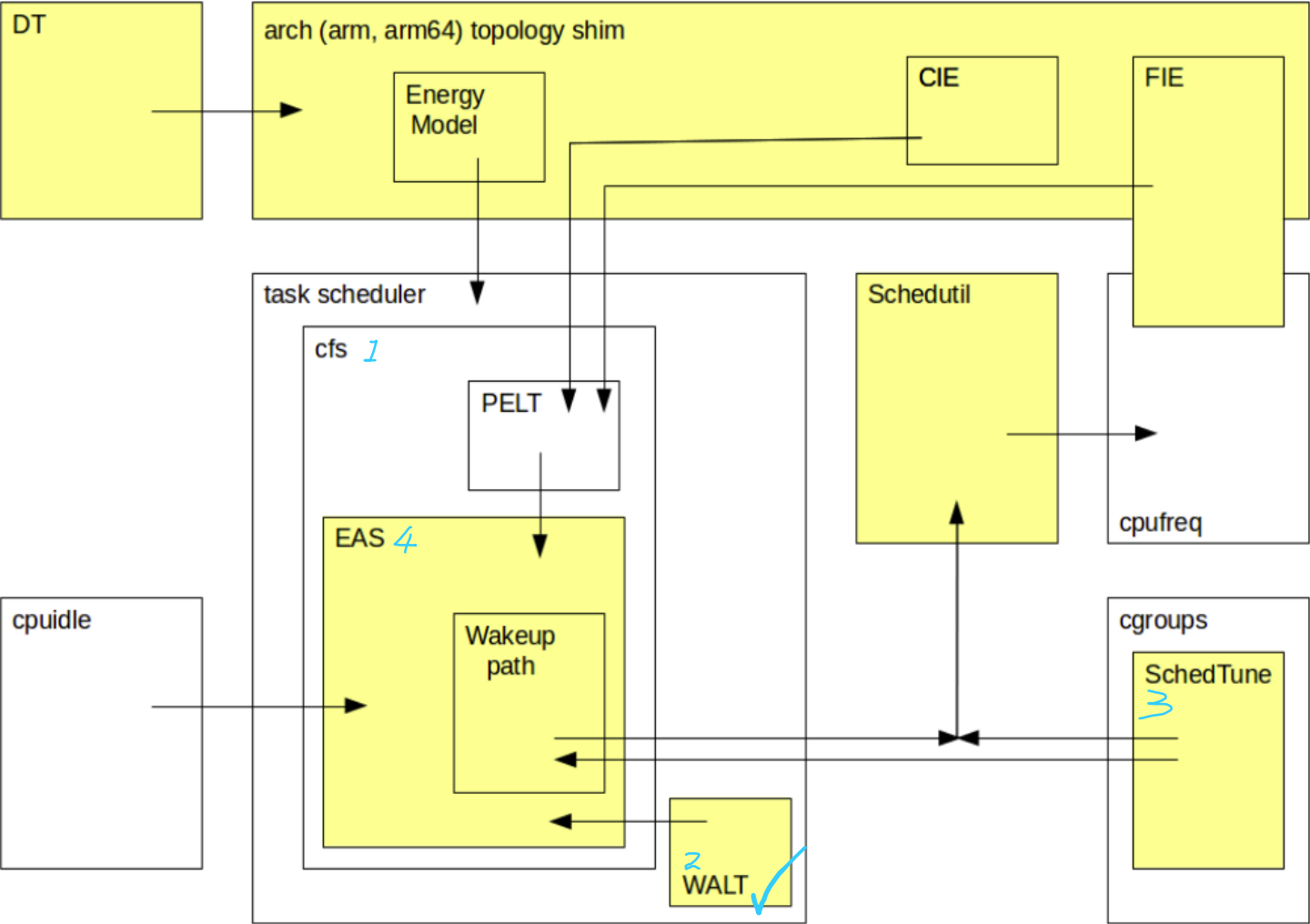


Figure 1 EAS building blocks in relation to Linux task scheduler, cgroups subsystem and related power management subsystems

Sched util - Scheduler driven DVFS

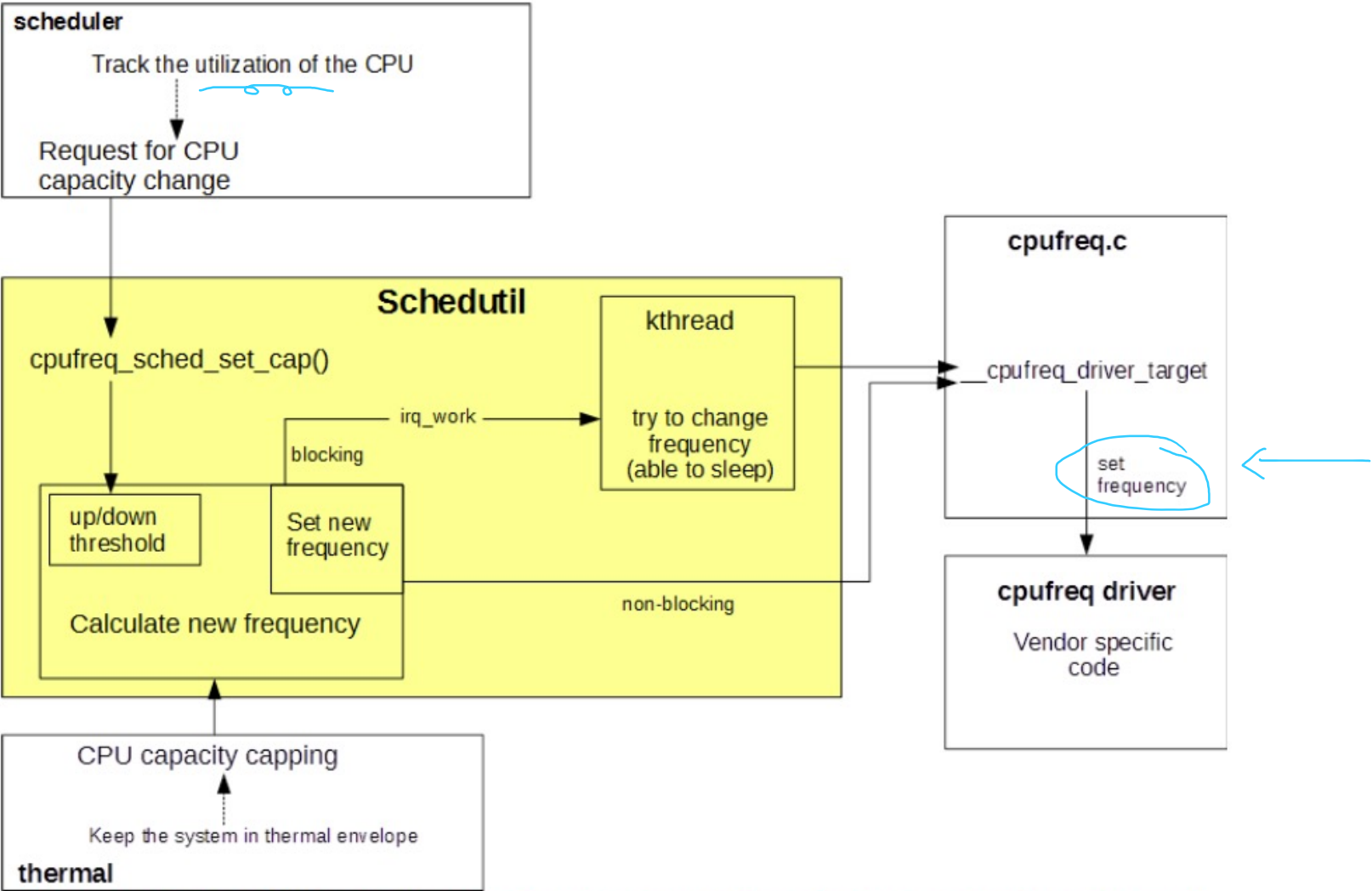


Figure 11 Schedutil block diagram showing connections between scheduler, thermal subsystem and existing CPUFreq

Sched Tune - Task classification and control

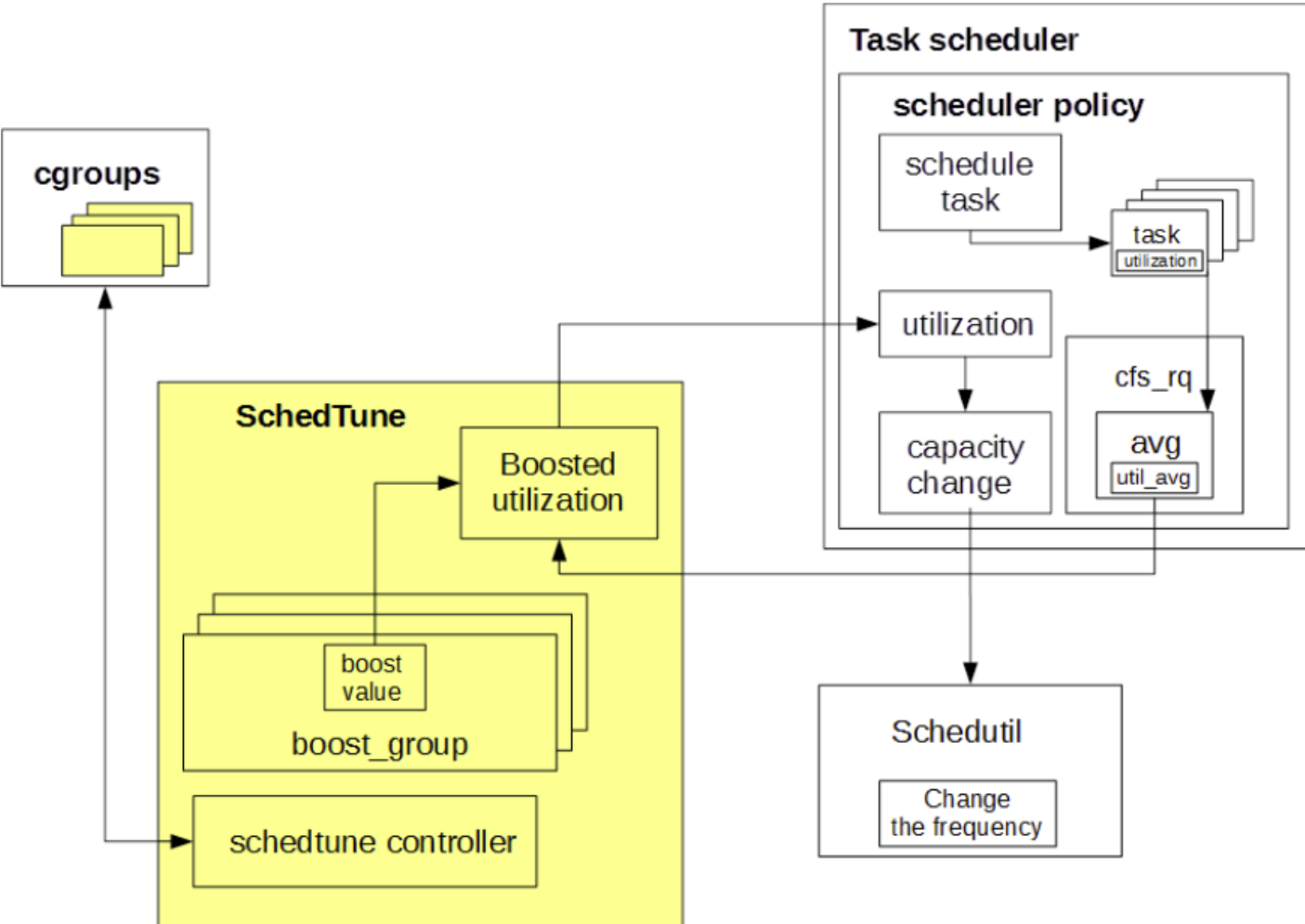


Figure 12 SchedTune block diagram showing components and connections between scheduler policy and boosted group's utilization concept.

This value is used to compute a margin to be added or removed to or from the utilization signal of a task/cpu. The value of the margin is calculated to provide a well-defined and expected user-space behaviour. For example, the following table reports the meaning of some specific boost values:

Boost value [%]	Meaning (e.g. run the task at a frequency corresponding to)
0	Minimum required capacity (max energy efficiency)
100	Maximum possible speed (min time to completion) (*)
50	Something in between the previous two configurations
-50	Half of the minimum required capacity
-100	Minimum available capacity (minimum OPP)

(*) minimum latency is not yet completely supported in the current ACK release, this feature is a work in progress and will be added in a following release.

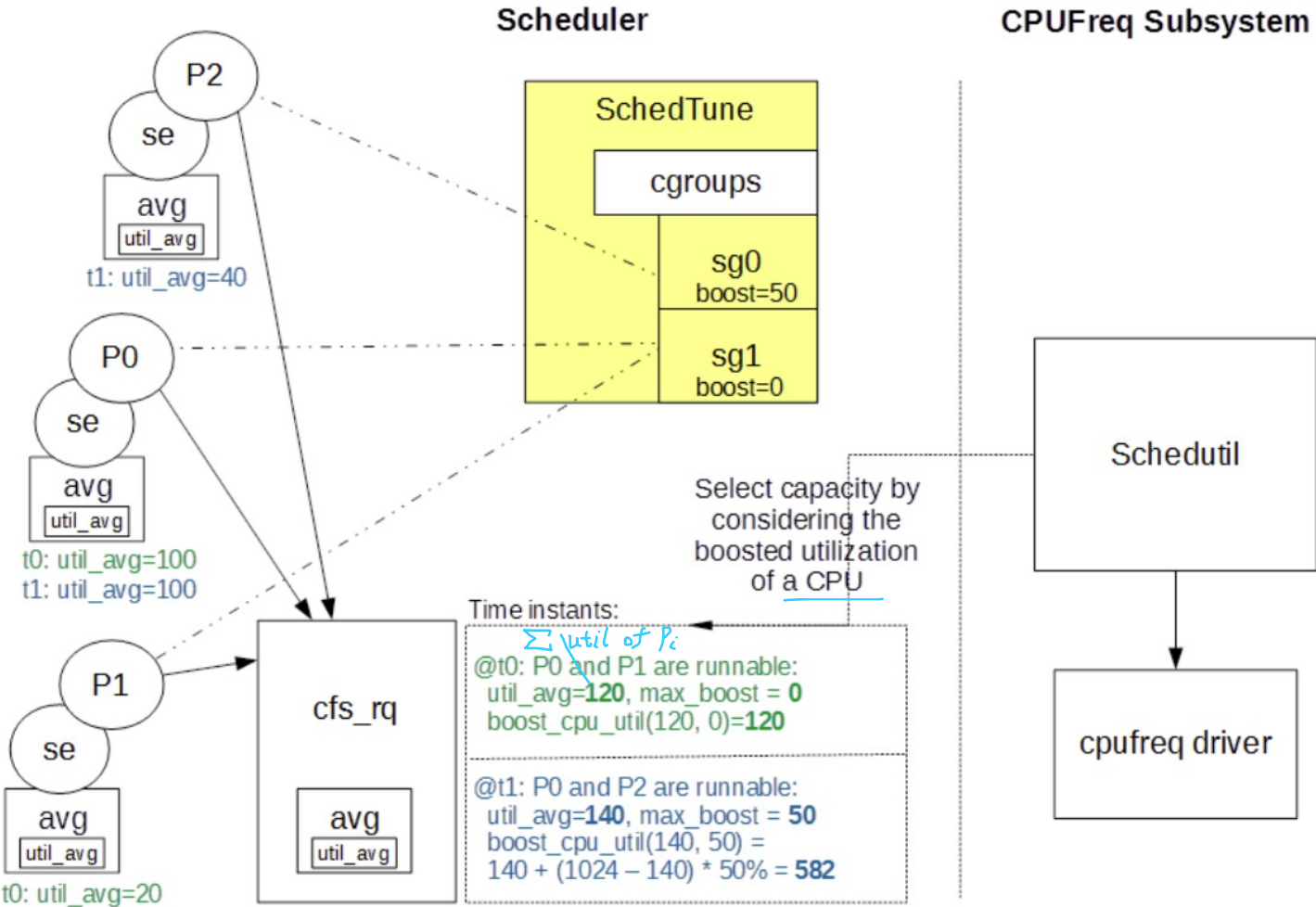


Figure 13 Flow diagram showing the state of the SchedTune and kernel scheduler in time: t0, t1.

Existing Linux CFS and EAS implementation

Aim: Energy-Efficiency and Performance.
Fairness

Terminology

- Energy = [joule] (resource like a battery on powered devices)
- Power = energy/time = [joule/second] = [watt]

The goal of EAS is to minimize energy, while still getting the job done.

maximize: *performance [inst/s] / Power [W]*

minimize: *Energy [J] / instruction*

EAS:

To break the tie between several good CPU candidates and pick the one that is predicted to yield the **best energy consumption** without harming the system's throughput when deciding which a task should run during wake-up.

Rely on specific elements about the platform's topology, the 'capacity' of CPUs, and their respective energy costs.

<https://docs.kernel.org/scheduler/sched-design-CFS.html>

<https://docs.kernel.org/scheduler/sched-energy.html>

Existing Linux CFS and EAS implementation

Qualcomm HMP scheduler

- Task demand D_{task} is the contribution of a task's running time to a window

1.25*D -- 80% capacity

- $D_{task} = \frac{\text{delta_time} \times \text{cur_freq}}{\text{max_possible_freq}}$
 - delta_time - time of task running on a core in a period of time
 - cur_freq - the current frequency of the core this task is running on
 - max_possible_freq is the maximum possible frequency across **all** cores

- Calculated over N sliding windows (N is a parameter)
 - E. g. the average demand $D_{avg} = (D_1 + \dots + D_N)/N$
 - The best result is achieved with $D = \max\{D_{avg}, D_1\}$

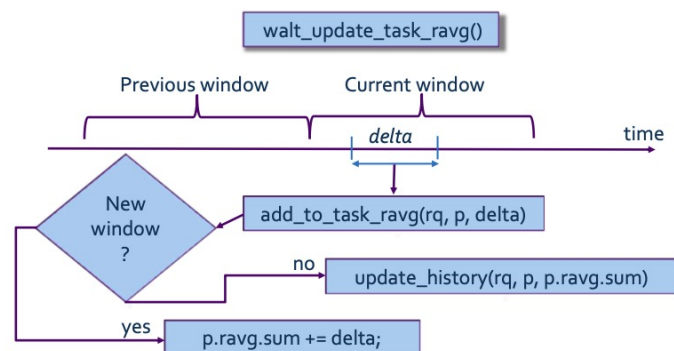
- We already account for difference in maximum frequency
 - D_{task} is calculated in regard to maximum frequency across all cores

- We also need to account for higher performance of big cores

- $D_{task, scaled} = D_{task} \cdot \frac{\text{rq} \rightarrow \text{efficiency}}{\text{max possible efficiency}}$
 - Efficiency* is a per-runqueue parameter
 - Usually big cores are considered 2x more effective

WALT: demand contribution calculation

Konsulko Group



Multimedia

To eliminate the above inefficiency of the conventional task scheduler, we propose an advanced task scheduler for heterogeneous mobile multi-core processors. **Our proposed task scheduler** allocates multimedia applications to the small cores and non-multimedia applications to the big cores.

Since recent smart devices have employed dedicated hardware decoders for multimedia applications

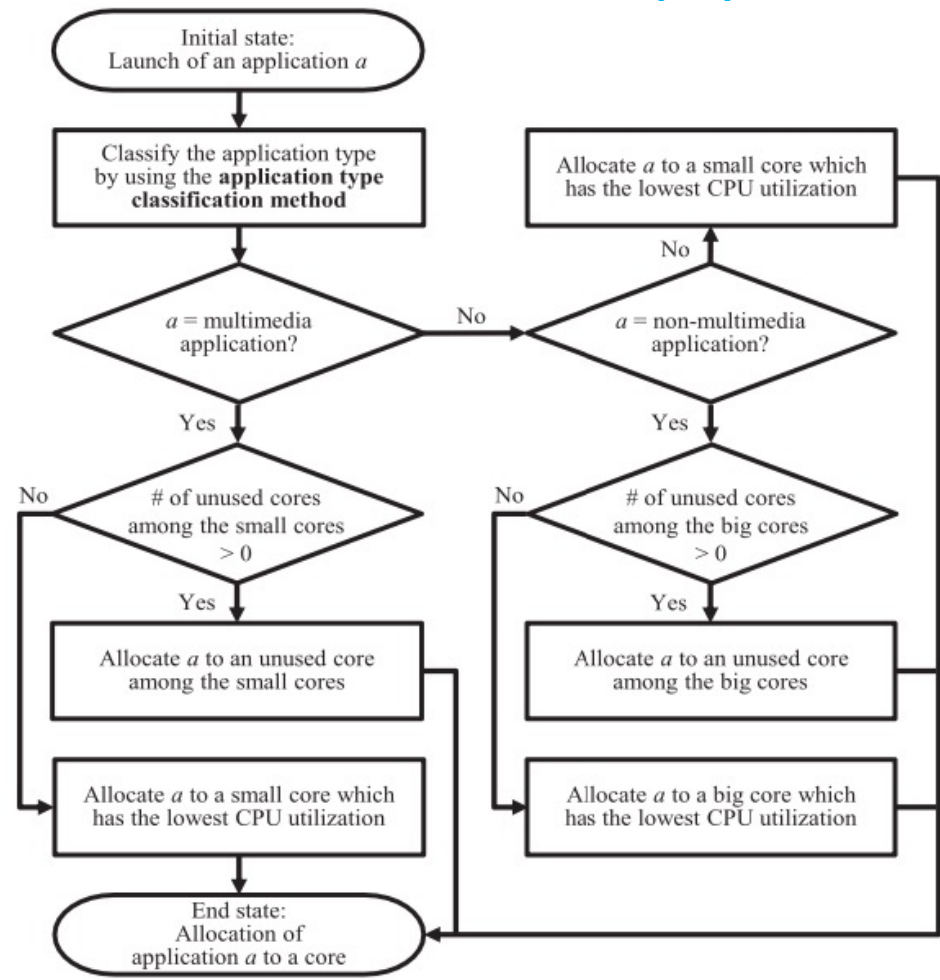


Fig. 5. Allocation algorithm of our proposed task scheduler.

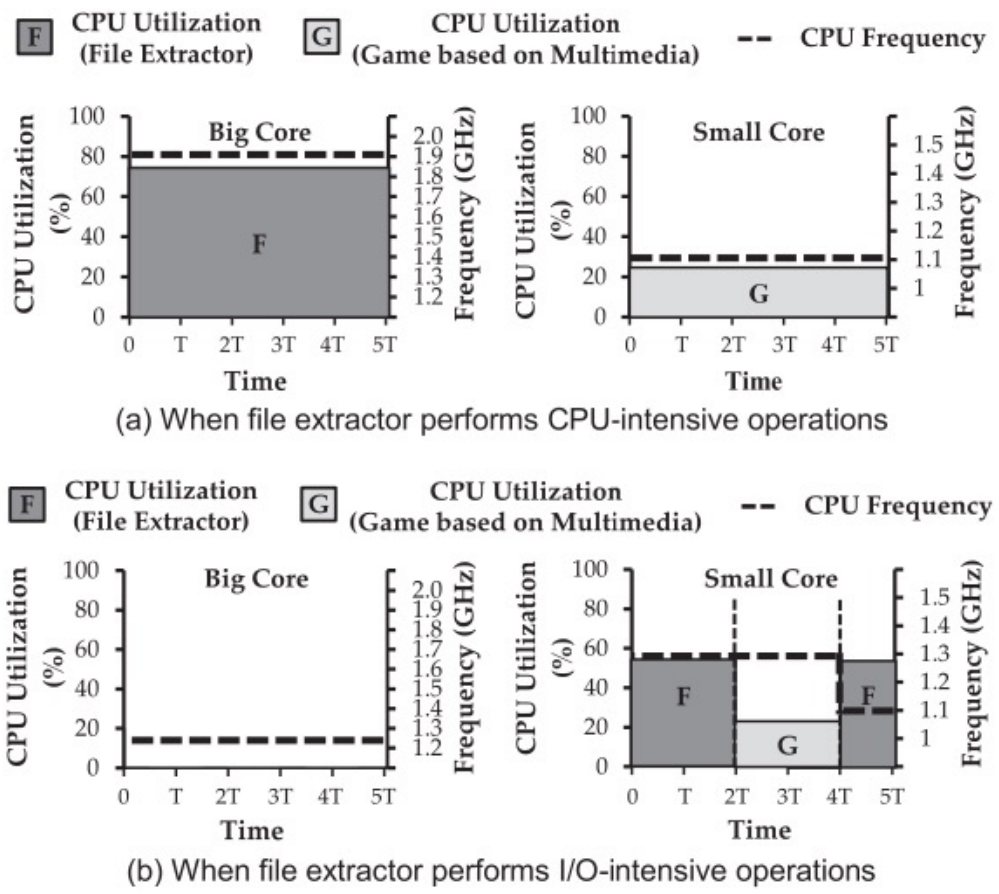


Fig. 2. Example of conventional task scheduler.

WAEAS

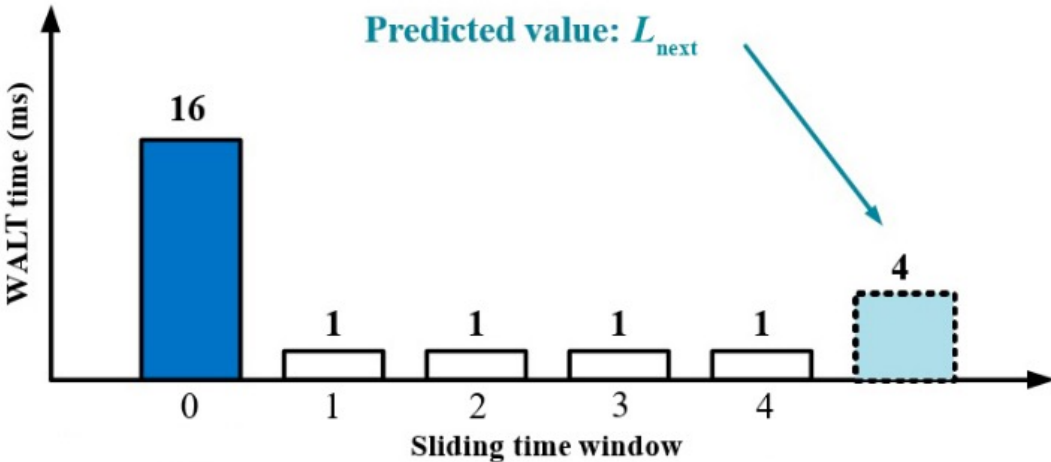
This paper analysed performance energy efficiency in the dynamic scheduling of latency-sensitive and non-latency-sensitive tasks for heterogeneous multicore wearable systems. By aiming to improve the energy efficiency of workload prediction, CPU core selection, and load balancing without affecting performance, the authors further optimized the EAS scheduler and proposed WAEAS, such as

BES-WALT **Basic exponential smoothing-based WALT algorithm,**

$$L = \text{time}_{\text{exec}} \times \frac{\text{freq}_{\text{curr}}}{\text{freq}_{\text{max}}} \times \frac{\text{IPC}_{\text{curr}}}{\text{IPC}_{\text{max}}}$$

L_t represents the predicted demand in the t -th sliding window,
 x_t represents the real demand in the t -th sliding window,
 α represents the smoothing factor, where α is in $[0,1]$

$$L_{t+1} = \alpha \times x_t + (1 - \alpha) \times L_t$$



$$\begin{aligned} L_{\text{recent}} &= 1 \\ L_{\text{avg}} &= (16 + 1 + 1 + 1 + 1) / 5 = 4 \\ L_{\text{next}} &= \max\{L_{\text{recent}}, L_{\text{avg}}\} = \max\{1, 4\} = 4 \end{aligned}$$

WAEAS

This paper analysed performance energy efficiency in the dynamic scheduling of **latency-sensitive and non-latency-sensitive tasks** for heterogeneous multicore wearable systems. By aiming to improve the energy efficiency of workload prediction, CPU core selection, and load balancing without affecting performance, the authors further optimized the EAS scheduler and proposed WAEAS, such as

- BES-WALT Basic exponential smoothing-based WALT algorithm,
- Energy-aware CPU selection algorithm,

Sched Tune

Latency-Sensitive Task: decrease searching speed to get minimum performance capacity and the lowest idle state

Non-Latency-Sensitive Task: not choose a lower utilization, but “task packing strategy” (on the little cores)

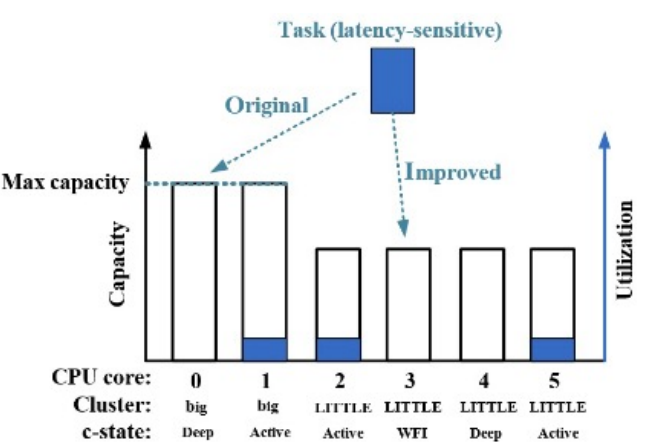


Fig. 5 Example of CPU selection for latency-sensitive tasks.

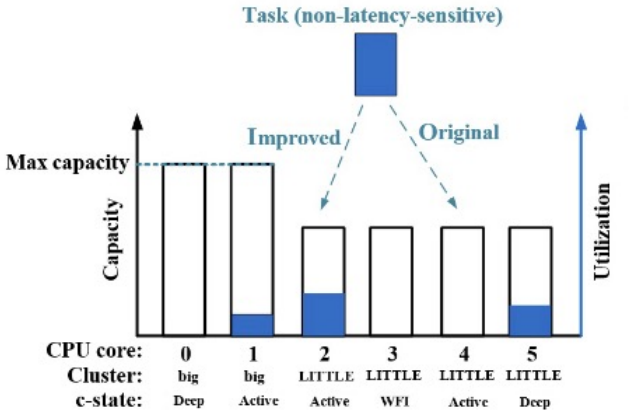


Fig. 6 Example of CPU selection strategy for non-latency-sensitive task.

WAEAS

This paper analysed performance energy efficiency in the dynamic scheduling of latency-sensitive and non-latency-sensitive tasks for heterogeneous multicore wearable systems. By aiming to improve the energy efficiency of workload prediction, CPU core selection, and load balancing without affecting performance, the authors further optimized the EAS scheduler and proposed WAEAS, such as

- BES-WALT Basic exponential smoothing-based WALT algorithm,
- Energy-aware CPU selection algorithm,
- Batch processing strategy,**

select idle sibling()

The proposed strategy records which core the task was running on as the awakener at the last time, and runs the tasks in the system centralized on fewer cores, which not only reduces the migration of tasks in the system, but also makes more cores in the system be in the idle state.

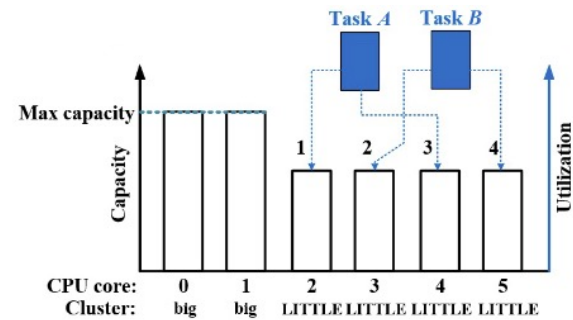


Fig. 7 Example of Tasks A and B awoken each other in the original algorithm: Task A (Core 2) wake up Task B (Core 3); Task B (Core 3) wake up Task A (Core 4); Task A (Core 4) wake up Task B (Core 5).

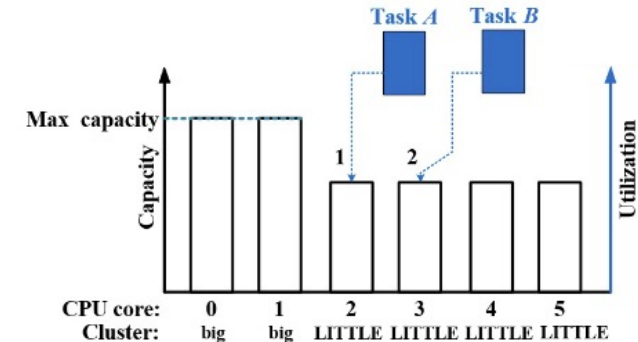


Fig. 8 Example of Tasks A and B waking up each other in the optimized algorithm: Task A (Core 2) wake up Task B (Core 3); Task B (Core 3) wake up Task A (Core 2); Task A (Core 2) wake up Task B (Core 3).

WAEAS

This paper analysed performance energy efficiency in the dynamic scheduling of latency-sensitive and non-latency-sensitive tasks for heterogeneous multicore wearable systems. By aiming to improve the energy efficiency of workload prediction, CPU core selection, and load balancing without affecting performance, the authors further optimized the EAS scheduler and proposed WAEAS, such as

BES-WALT Basic exponential smoothing-based WALT algorithm,

Energy-aware CPU selection algorithm,

Batch processing strategy,

Cluster-based load balancing (overutilized)

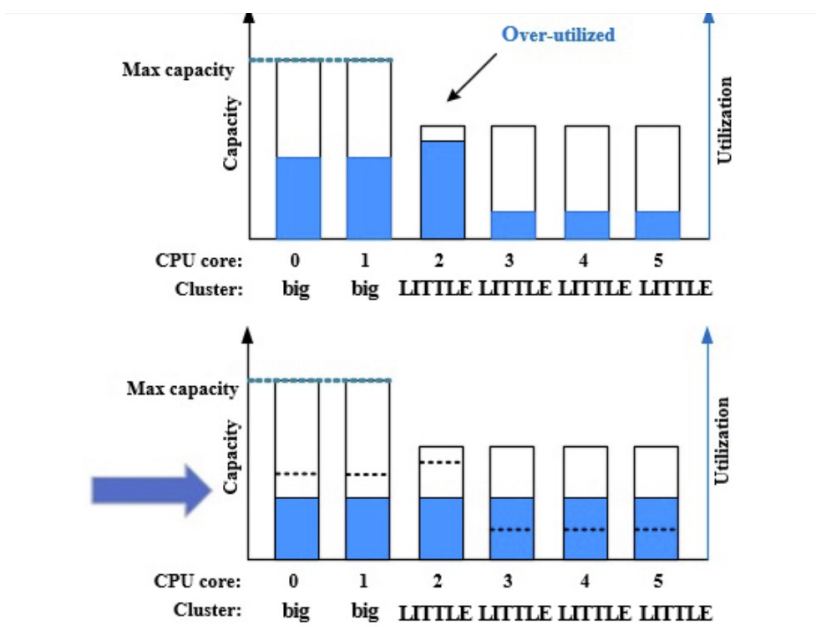


Fig. 9 Example of the original load balancing in EAS scheduler.

reduces the meaningless task migration by *local* load balancing sets an overutilized indicator in the **big** core scheduling subdomain and the **little** core scheduling subdomain, **respectively**, thereby indicating that whether the domain is overutilized or not. When a scheduling domain is overutilized, the current domain state is set to be overutilized first, and the load balancing in the current domain is given a priority

aperiodic high-load applications

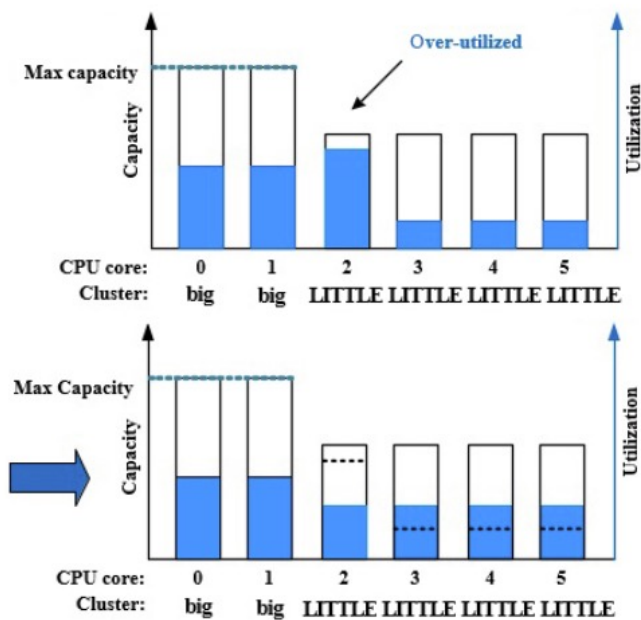


Fig. 10 Example of the cluster-based load balancing strategy.

Learning EAS

The Learning EAS adjusts the TARGET_LOAD used to set the CPU frequency and the sched_migration_cost used as the task migration criteria according to the characteristics of the running task through the **policy gradient reinforcement learning**.

Effects:

In LG G8 ThinQ, Learning EAS improved power consumption by 2.3% - 5.7%, hackbench results for process scheduling performance by 2.8% - 25.5%, applications entry time by 4.4% - 6.1%, and applications entry time under high CPU workload by 9.6% - 12.5%, respectively compared with EAS. This paper also showed that the Learning EAS is scalable by applying the Learning EAS to high- end and low-end chipset platforms of Qualcomm.Inc and MediaTek.Inc and improving power consumption by 2.8% - 7.8%, application entry time by 2.2% - 7.2%, respectively compared with EAS. Finally, this paper showed that the performance of CPU scheduling is improved gradually by the repetition of reinforcement learning.

Learning EAS

Motivation:

EAS generally provides good performance because many of the values which are basis for the operation are **fixed**.

In general, each task has its own characteristics for changing *workload*.

It also has unique properties for *transition between running state and sleep state*.

It is difficult to optimize the performance by using simple algorithm which changes the default value for the EAS operation when a task's *workload* or the *ratio of sleep and running states* gets to change above or below a certain value.

Solution:

Set the values of the EAS operation criteria dynamically according to the characteristics of the running tasks to improve performance, which can be achieved through **policy gradient reinforcement learning**.

$$\text{updating parameter vector} \quad \theta_{t+1} = \theta_t + \alpha \frac{\partial \rho}{\partial \theta_t} \quad \begin{matrix} \text{learning rate} & \text{policy performance measure} \end{matrix}$$

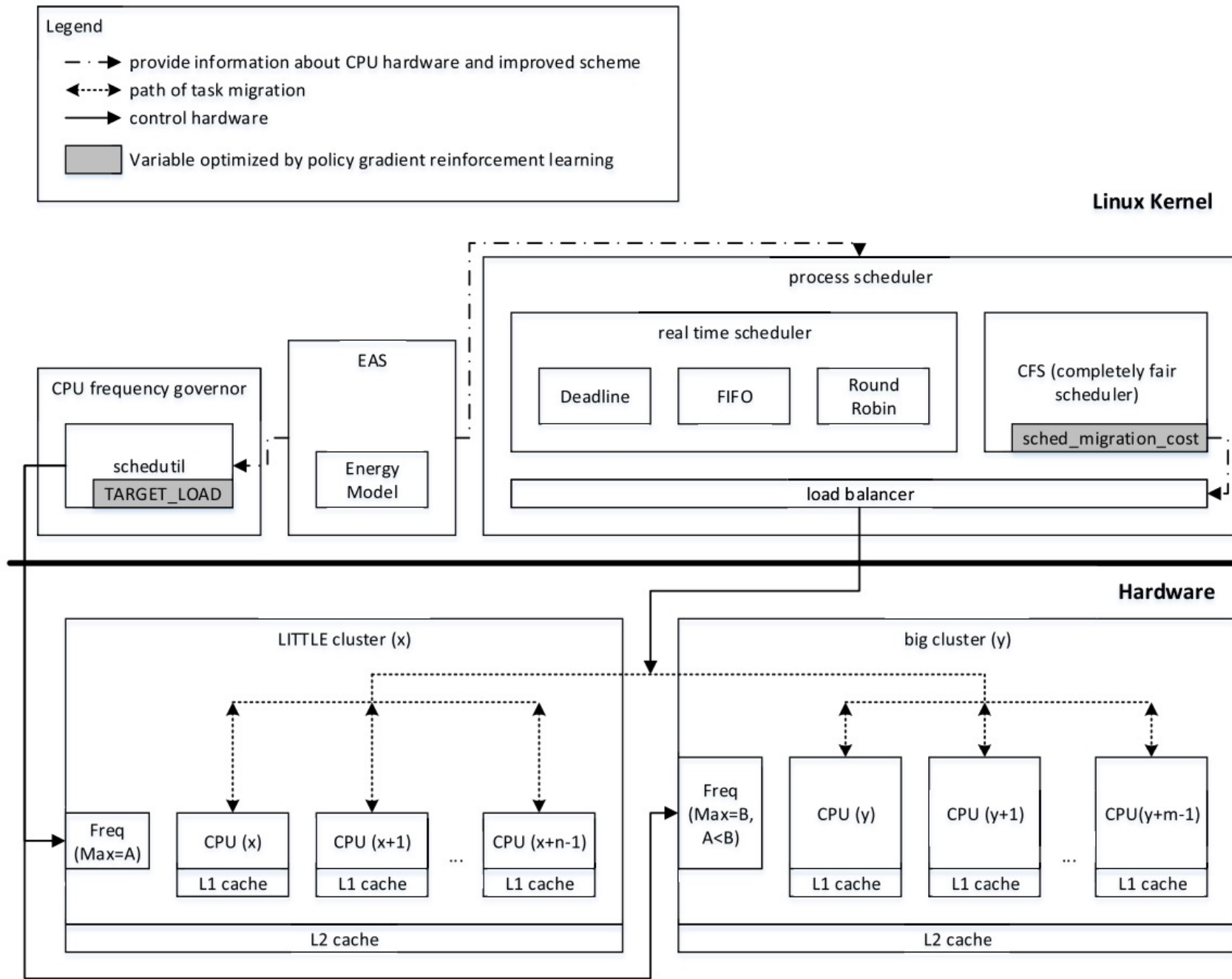


FIGURE 3. Block diagram of EAS and big.LITTLE chipset platform.

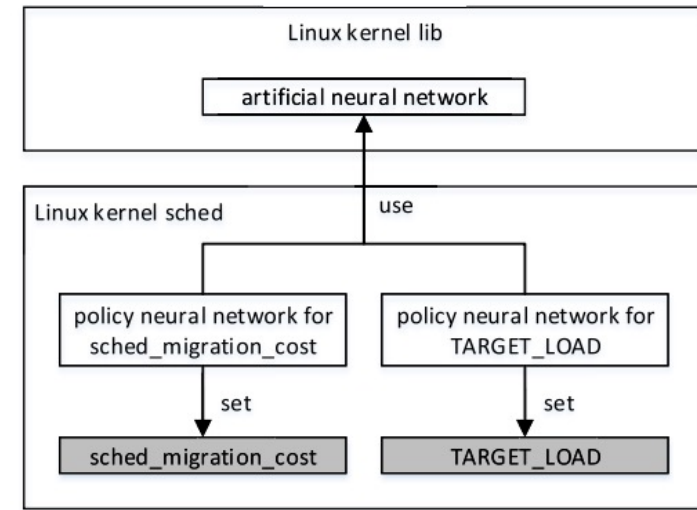


FIGURE 7. Diagram of policy gradient reinforcement learning for learning EAS.

cy Gradient Reinforcement Learning

Learning EAS

$$s_0 \xrightarrow{a_0} s_1 \xrightarrow{a_1} s_2 \xrightarrow{a_2} s_3 \xrightarrow{a_3} \dots$$

State s each CPU util

Action a increase or decrease TARGET_LOAD

Policy $\pi : S \times A \rightarrow [0,1]$ $\pi(a, s) = \Pr(a_t = a \mid s_t = s)$

the agent's action selection, π mapping from the states to the actions.

Stationary Probability $d^\pi(s) = \lim_{t \rightarrow \infty} P(s_t = s \mid s_0, \pi_\theta)$

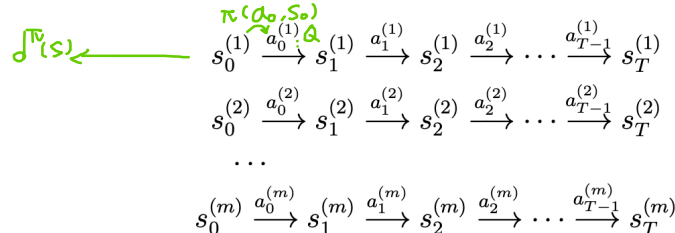
State-value function $V : S \times \pi \rightarrow \mathbb{R}$ $V^\pi(s) = E[R \mid s, \pi]$.

Action-value function $f_w/Q : S \times \pi \rightarrow \mathbb{R}$ Expected return

Performance Function ρ $\rho(\theta) = \rho^{\pi_\theta}$, θ space \rightarrow policies π_θ space

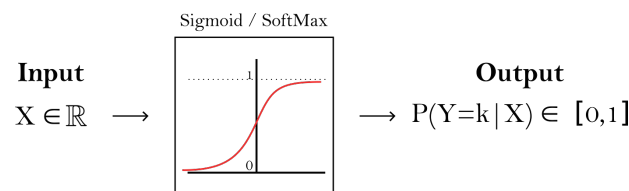
parameter θ the policy that *increases or decreases* TARGET_LOAD

which affects the next CPU frequency change.



Action Preference Function $h(s, a, \theta)$

$$\pi(a \mid s, \theta) = \frac{e^{h(s, a, \theta)}}{\sum_b e^{h(s, b, \theta)}}$$



Optimal Value Function

$$\min_{T_j} \sum_{j=1}^{N_{cpu}} P(T_j) \quad P(T_j) \approx a \cdot (T_j)^2 + b$$

$$\text{subject to } \sum_{i=1}^{N_{task}} W_i \leq \sum_{j=1}^{N_{cpu}} T_j \quad (6)$$

CPU workload \leq CPU throughput

In schedutil, each CPU workload is **Max** expressed as $C_{util} = \left(\frac{C_{time}}{C_{time} + I_{time}} \right) \cdot MC_{util}$

largest CPU util in the cluster

$$\text{Next } NC_{freq} \cdot (0.8) = \left(\frac{LC_{util}}{MC_{util}} \right) \cdot MC_{freq}$$

$$NC_{freq} = \left(\frac{LC_{util}}{MC_{util}} \right) \cdot MC_{freq} \cdot (1.25)$$

Reward function $J(\theta) = \sum_{s \in S} d^\pi(s) V^\pi(s) = \sum_{s \in S} d^\pi(s) \sum_{a \in A} \pi_\theta(a \mid s) Q^\pi(s, a)$

Policy Gradient Theorem

$$\frac{\partial \rho}{\partial \theta} = \sum_s d^\pi(s) \sum_a \frac{\partial \pi(s, a)}{\partial \theta} f_w(s, a) \quad (11)$$

d^π is the stationary distribution of states under π , which is independent of s_0 for all policies, and f_w is approximation of quality function. State s is each CPU util and a is the action to increase or decrease TARGET_LOAD. There are two typical

current TARGET_LOAD change value

$$\frac{\partial \rho}{\partial \theta_t} = tl_t(\text{freq}_t - \text{freq}_{t-1})$$

Direct policy search: $\theta_{t+1} = \theta_t + \alpha \frac{\partial \rho}{\partial \theta_t}$
Updating parameter vector α learning rate

Thus when $\pi(s, a)$ is **partially differentiated against θ** , it is the rate of the CPU **frequency change**.

Learning EAS

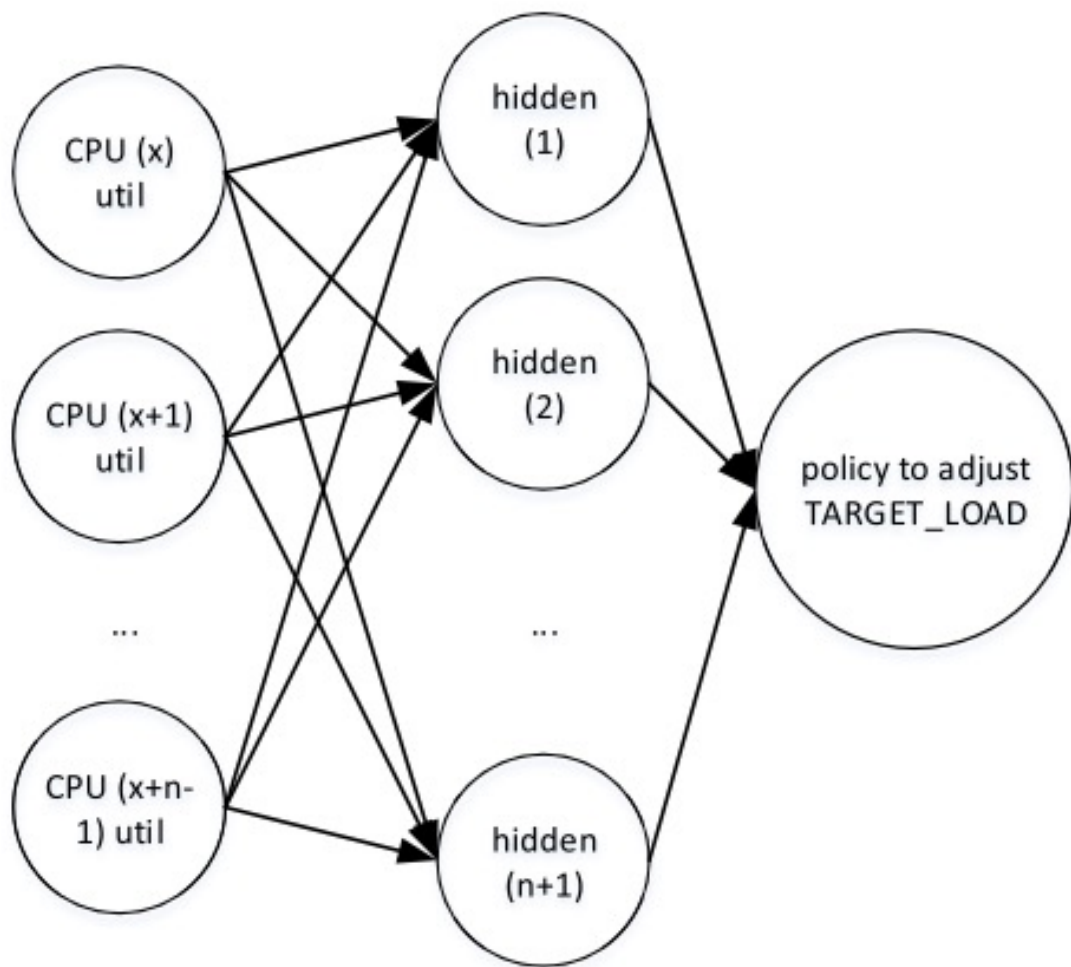


FIGURE 4. Configuration of policy neural network to set TARGET_LOAD for cluster (x).

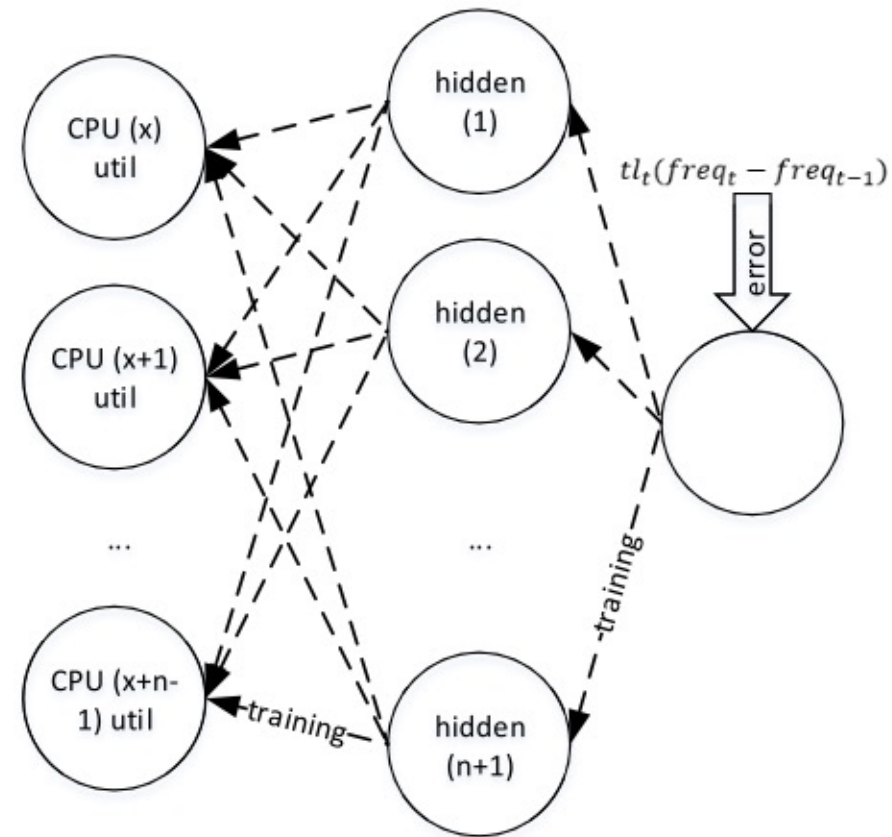


FIGURE 5. Training of policy neural network to set TARGET_LOAD for cluster (x).

$$\frac{\partial \rho}{\partial \theta_t} = tl_t(\underline{freq_t - freq_{t-1}})$$

Learning EAS

Algorithm 1 Pseudocode to set TARGET_LOAD for Cluster (x) Using Policy Gradient Reinforcement Learning

```
Initialize policy neural network
CPU_utils[n] = {0, 0, ..., 0}
prev_CPU_freq = cur_CPU_freq = 0
repeat
    i = 0
    while i < n do
        CPU_utils[i] = Get ith CPU load
        ++ i
    end
    result = Run policy neural network(CPU_utils)
    if (threshold(12000) < result) then
        tl = -0.0023
    else then
        tl = 0.0023
    end
    TARGET_LOAD += tl
    prev_CPU_freq = cur_CPU_freq
    Sleep 100ms
    Get cur_CPU_freq
    error = tl * (cur_CPU_freq - prev_CPU_freq)
    Training policy neural network by back
    propagation(error)
until system has been shutdown
```

Algorithm 2 Pseudocode to set Sched_Migration_Cost Using Policy Gradient Reinforcement Learning

```
Initialize policy neural network
CPU_load_vari[m] = {0, 0, ..., 0}
prev_total_CPU_load_vari = total_CPU_load_vari = 0
mc = 0
prev_mig_count = mig_coun = 0
prev_mig_success = mig_success = 0
repeat
    Sleep 7000ms
    total_CPU_load_vari = 0
    i = 0
    while i < m do
        CPU_load_vari[i] = Get the variance of CPU
        workloads in ith cluster
        total_CPU_load_vari += CPU_load_vari[i]
        ++ i
    end
    mig_count = Get task migration attempts counts
    mig_success = Get task migration success counts
    cur_error = (total_CPU_load_vari /
    (mig_count + mig_success))
    prev_error = (prev_total_CPU_load_vari /
    (prev_mig_count + prev_mig_success))
    error = -mc * (cur_error - prev_error)
    Training policy neural network by back
    propagation(error)
    result = Run policy neural network
    (mig_count, mig_success, CPU_load_variance)
    if (threshold(160000) < result) then
        mc = -17027
    else then
        mc = 17027
    end
    sched_migration_cost += mc
    prev_total_CPU_load_vari = total_CPU_load_vari
    prev_mig_count = mig_count
    prev_success_count = success_count
until system has been shutdown
```

Thank you!
Best of Luck~

Peter Hu

Huawei Technologies Research & Development

CPU Team

