# Linux Power

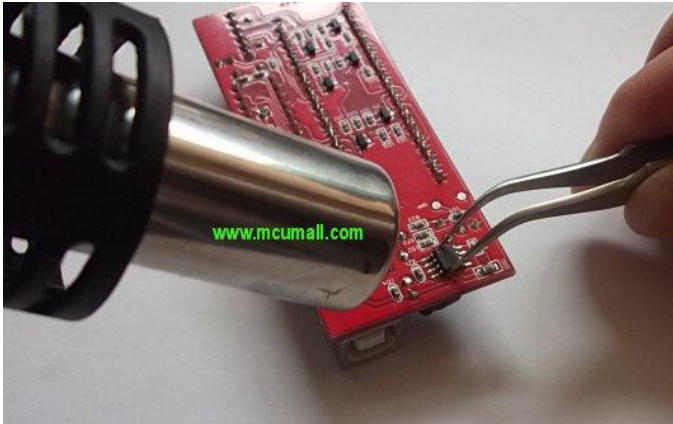## Basic, Linux PM, DVFS, OPP, Runtime PM, PM QoS
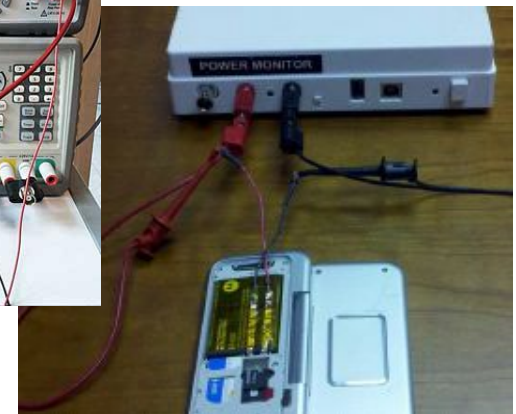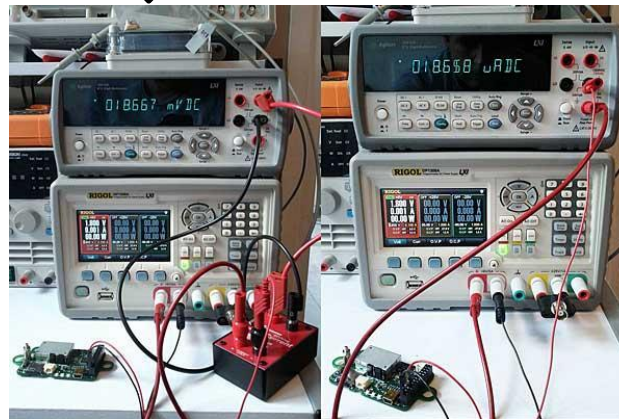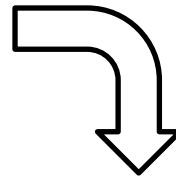
국민대학교 임베디드 연구실
경 주 현

# Outline

- **Processor Power Management**

- **Linux Power Management**

# How to reduce power?

# How to reduce power?



WIFI, Modem chip, LCD, Touch ....



- Power Rail Check
- Total Current Check

# Background

- **Energy**
  - joule

- **Power**
  - watt

- **Processor Power Management**
  - P-state, C-state

# Processor Power Management

# Silicon Power Management Features



- **P – States are a sub-state of the C0 state**
  - They offer reduced power consumption while the processor is executing code.

- **C-States are a sub-state of the S0 state**
  - They offer reduced power consumption while the system is fully on.

- **S-States**
  - Suspend to idle, Standby, Suspend to ram, suspend to disk

# Processor Power Management

- **CPU C-State (Idle Power State)**

  – C0  - Full On State

  – C1 – Auto Halt State

  – C2, C4 (C3, C4, C5), C6


- **CPU P-State (Performance State)**

| | Active state | Sleep states | | | |
|---|---|---|---|---|---|
| | C0 | C1/C1E | C3 | C6 | C7 |
| | Operating | Halt | Sleep | Deep Sleep | |
| Core clock | ⊓⊔ | off | off | off | off |
| PLL | ⊓⊔ | ⊓⊔ | off | off | off |
| Core caches | | | flushed | flushed | flushed |
| Shared cache | | | | | flushed |
| Wakeup time* | active | | | | |
| Core Idle power* | | | | ~ 0 | < C6 |

\* Rough approximation

http://download.intel.com/embedded/technology/PowerManagementforEmbeddedApps.pdf

# Silicon Power Management Features

- **Advanced Configuration and Power Interface (ACPI)**

Performance state(P-states)
P0 : highest performance, highest power
P1 ~ Pn : lower frequency and voltages

Throttling states(T-states)
T1~Tn : change only frequency not voltage
for thermal

```
┌─────────────────────────────┐
│  Performance  │  C0  │ Throttling │
│  State Px     │      │            │
└─────────────────────────────┘
        C2
   C1          C3
```

Idle(stops CPU main internal clocks via HW)

Idle(stops CPU main internal clocks via SW)

Idle(stops All CPU internal clocks, not kept cache coherent)

# Linux Power Management

# Tree view - Dynamic and Static PM

KESL Kookmin Univ. Embedded System Lab

# Linux Driver Model

- **Example**

```
static struct device_driver eepro100_driver = {
        .name           = "eepro100",
        .bus            = &pci_bus_type,

        .probe          = eepro100_probe,
        .remove         = eepro100_remove,
        .suspend        = eepro100_suspend,
        .resume         = eepro100_resume,
};
```

# Linux Power Framework



https://wiki.linaro.org/WorkingGroups/PowerManagement/Doc/Architecture

# Linux Driver Model



https://wiki.linaro.org/WorkingGroups/PowerManagement/Doc/Architecture

# Tree view - Dynamic and Static PM

KESL Kookmin Univ.
Embedded System Lab

# General System Suspend/Resume Control Flow



Power Management In The Linux* Kernel - Rafael J. Wysocki

# Problem

- **How to reduce system idle power?**

# Problem

- **How to reduce system idle power?**

- **What about CPU Hotplugging?**

KESL Kookmin Univ.
Embedded System Lab

# Problem

- **How to reduce system idle power?**

- **What about CPU Hotplugging?**
  - **Expensive latency**

- **Solution?**

# Problem

- **How to reduce system idle power?**

- **What about CPU Hotplugging?**
    - **Expensive latency**

- **Solution?**

    **S2I(Suspend to Idle)**

# Suspend-to-Idle



Suspend

**Call Notifiers**

**Freeze Tasks**

**Device Suspend**

.prepare()

.suspend()

.suspend_late()

Interrupt handlers run
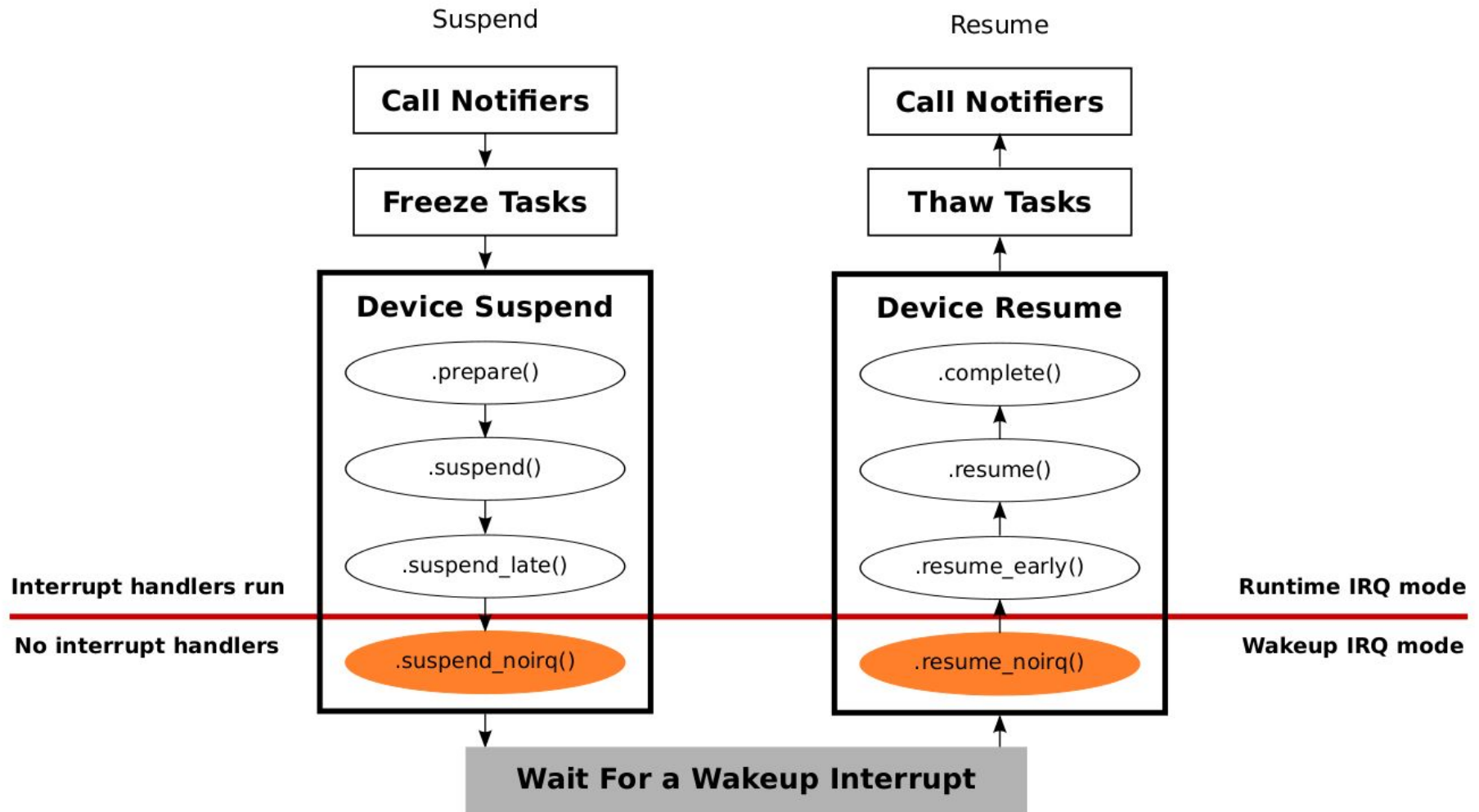
No interrupt handlers

.suspend_noirq()

Resume

**Call Notifiers**

**Thaw Tasks**

**Device Resume**

.complete()

.resume()

.resume_early()

Runtime IRQ mode

Wakeup IRQ mode

.resume_noirq()

**Wait For a Wakeup Interrupt**

https://www.linaro.org/blog/suspend-to-idle/

KESL Kookmin Univ. Embedded System Lab
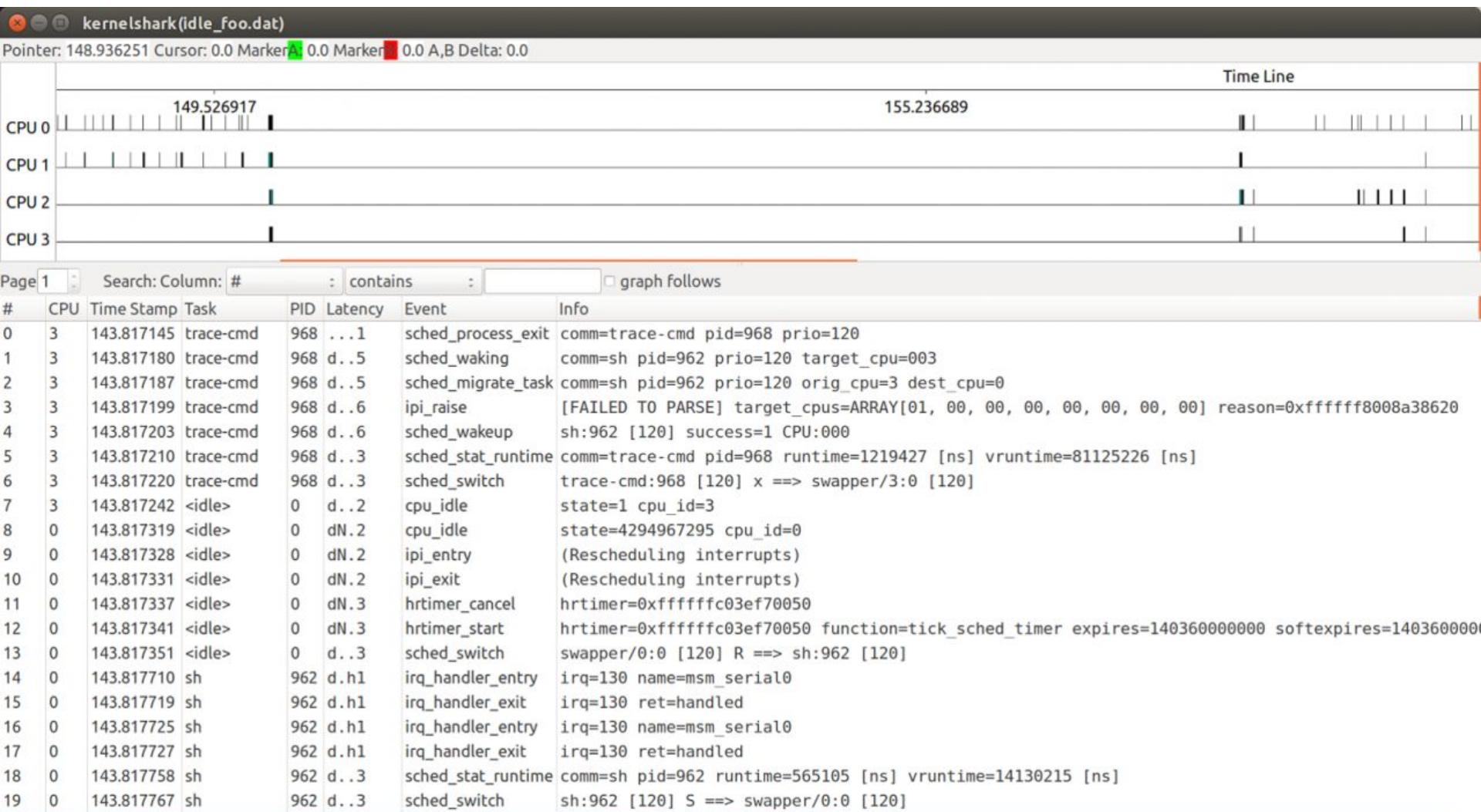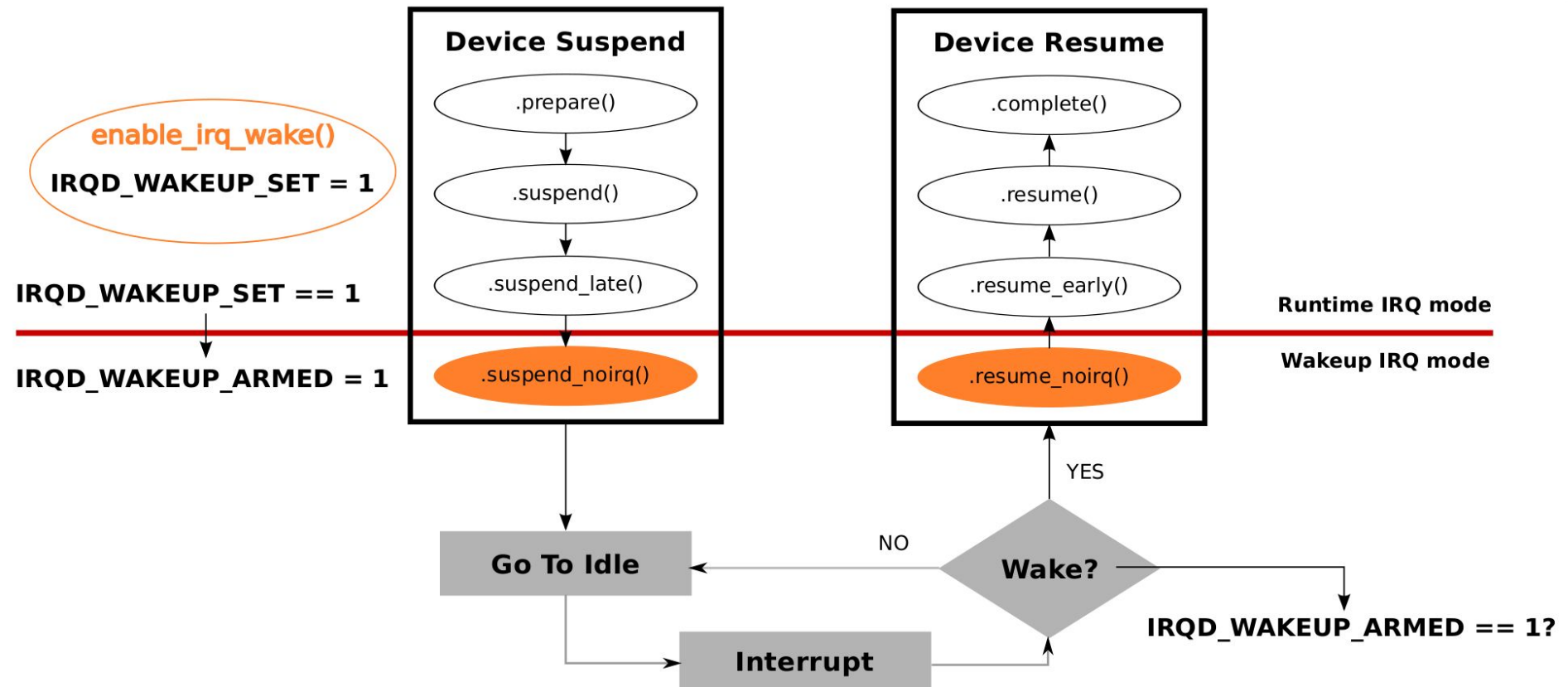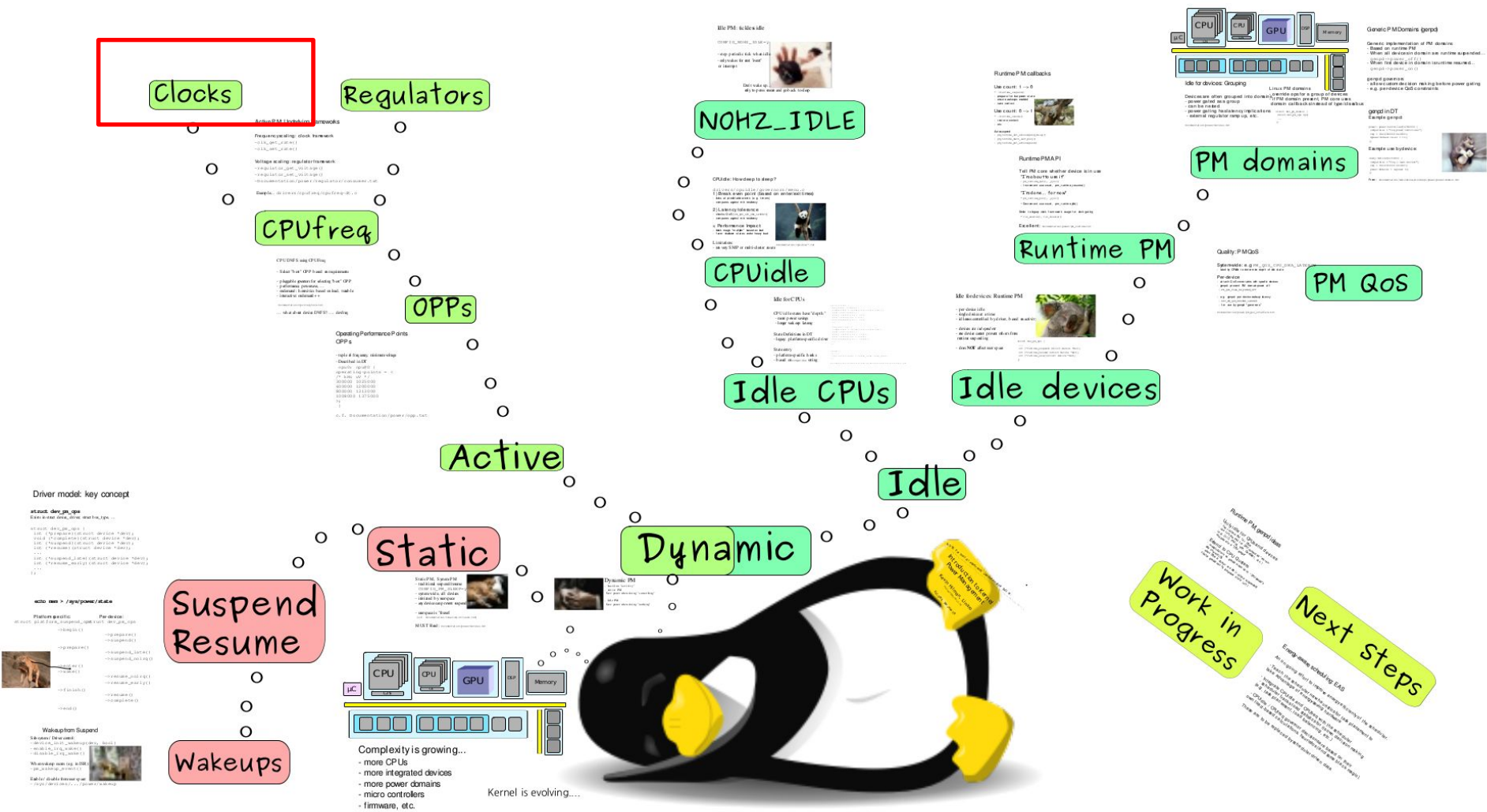
# Suspend-to-Idle

# Suspend-to-Idle

# Suspend-to-Idle Wakeup



"What Is Suspend-to-Idle and How To Make It Work", Rafael J. Wysocki

# Tree view - Dynamic and Static PM

# Linux Power Framework



https://wiki.linaro.org/WorkingGroups/PowerManagement/Doc/Architecture

# Clock Framework

- **Clock Domains**

# Clock Framework – function list

- **clk_get**
  - Lookup and obtain a reference to a clock producer

- **clk_enable**
  - Inform the system when the clock source should be running.

- **clk_disable**
  - Inform the system when the clock source is no longer required.

- **clk_put**
  - "free" the clock source

# Tree view - Dynamic and Static PM

# Linux Power Framework



https://wiki.linaro.org/WorkingGroups/PowerManagement/Doc/Architecture

# Regulator Framework - 1

# Regulator Framework - 4



- Regulators have power constraints to protect hardware
- Sysfs interface (/sys/class/regulator/.../)
- Regulator independent abstraction (drivers/regulators/core.c)
- **regulator_get** - lookup and obtain a reference to a regulator
- **regulator_put** - "free" the regulator source
- **regulator_enable** - enable regulator output
- **regulator_disable** - disable regulator output

# Regulator Framework - 2



Regulator-1 supplies Domain 1
Regulator-2 supplies Domain 2
Regulator-3 supplies Domain 3

# Tree view - Dynamic and Static PM

KESL Kookmin Univ. Embedded System Lab

# Linux Power Framework



https://wiki.linaro.org/WorkingGroups/PowerManagement/Doc/Architecture

# OPP - Operating Performance Point

- **Tuple of frequency, minimum voltage.**
- **Described in DT**

```
cpu0: cpu@0 {
operating-points = <
/* kHz uV */
300000 1025000
600000 1200000
800000 1313000
1008000 1375000
>;
  }
```

# Linux Power Framework

# Tree view - Dynamic and Static PM

KESL Kookmin Univ. Embedded System Lab

# CPU Frequency Governor
# - 4.5 and Earlier version of Linux -

# CPU Frequency Subsystem

- **Linux CPU frequency subsystem manages CPU frequency scaling.**

# CPU Frequency Subsystem

- **Comprise three parts**
  - Governor
    - Implemented CPU frequency scaling policy
    - Governors in Linux Kernel
      - Performance, Powersave, Userspace, Ondemand, Conservative, Interactive
  - Core
    - Provide generic interface for governor and platform driver to communicate with each other

  - Platform driver
    - SoC-chip specific cpu frequency scaling driver
      - Thermal throttle, DVFS, and so on

KESL Kookmin Univ.
Embedded System Lab

# Governors In Linux Kernel

- **Powersave governor**
  - Set frequency statically to lowest available frequency.

- **Performance governor**
  - Set frequency statically to highest available frequency.

- **Userspace governor**
  - Set frequency manually.

- **Ondemand governor**
  - Scale frequency based on CPU utilization.

- **Interactive governor**
  - More responsive to interactive workloads.

# Governors In Linux Kernel



"CPUFreq and The Scheduler Revolution in CPU Power Management", Rafael J. Wysocki

# DVFS in Linux (cpufreq)

# Ex) Interactive Governor

- **<u>Interactive</u>**



**Interactive**

# Interface



CPU-2 로드 증가

CPU-2 Frequency 증가

# Boost

# Floor validate_time



floor_validate_time
Default : 80ms

# Deferrable timer



하늘색
CPU 6 CPU LOAD

IDLE 진입, 하지만
CPU Load는 높다.

# Deferrable timer



분홍색
CPU 4 CPU LOAD

IDLE 진입, 하지만
CPU Load는 높다.

# DVFS in Linux (cpufreq)

- **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



"EAS Update", Amit Kucheria et. ,al. ARM

# Tree view - Dynamic and Static PM

KESL Kookmin Univ. Embedded System Lab

# Linux Power Framework



https://wiki.linaro.org/WorkingGroups/PowerManagement/Doc/Architecture

# PM_QOS

- **Goal : to reduce energy**

- **A coordination mechanism**
  - User's performance needs

- **Examples**
  - PM QoS with cpuidle
  - Per-device PM QoS with runtime pm.

# Linux Power Framework



https://wiki.linaro.org/WorkingGroups/PowerManagement/Doc/Architecture

# Tree view - Dynamic and Static PM

KESL Kookmin Univ. Embedded System Lab

# Processor Power Management



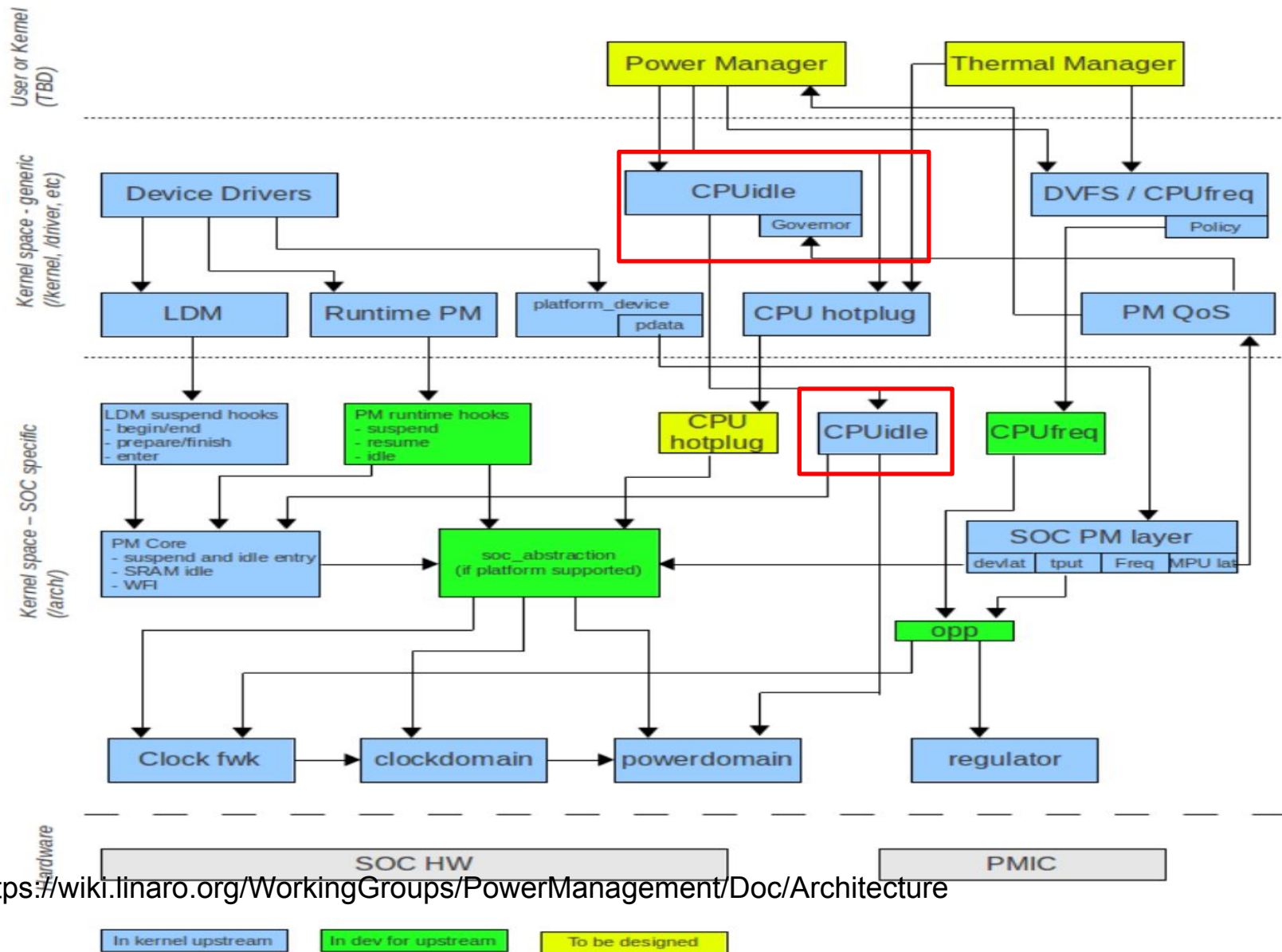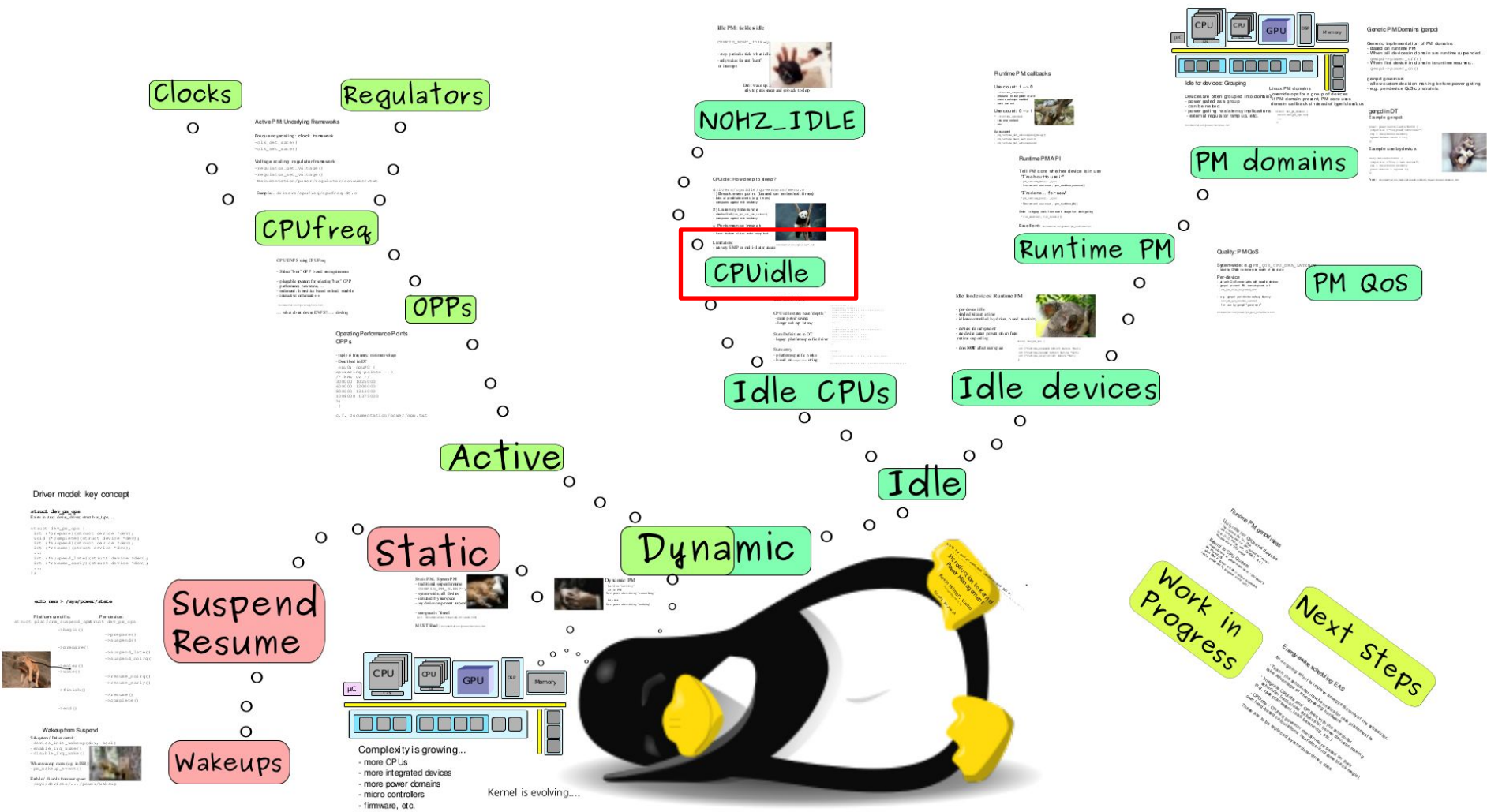| | Active state | Sleep states | | | |
|---|---|---|---|---|---|
| | C0 | C1/C1E | C3 | C6 | C7 |
| | Operating | Halt | Sleep | Deep Sleep | |
| Core clock | ⊓⊔ | off | off | off | off |
| PLL | ⊓⊔ | ⊓⊔ | off | off | off |
| Core caches | (active) | (active) | flushed | flushed | flushed |
| Shared cache | (active) | (active) | (active) | (active) | flushed |
| Wakeup time* | active | (small) | (medium) | (large) | (largest) |
| Core Idle power* | (full) | (reduced) | (low) | ~ 0 | < C6 |

* Rough approximation

KESL Kookmin Univ. Embedded System Lab

# CPU idle - C-state management

# CPU idle

# Menu governor: How deep to sleep?

- **Energy break**
  - next timer event
  - historic behavior

- **Performance impact**
  - nr_iowaiters, cpu_load

- **Latency    tolerance**
  - checks    QoS    (PM_QOS_CPU_DMA_LATENCY)

KESL Kookmin Univ.
Embedded System Lab

# Per-device Power Management

# Tree view - Dynamic and Static PM

KESL Kookmin Univ. Embedded System Lab

# Linux Power Framework

# Runtime Power Management

- Why's Runtime PM

- What's Runtime PM

- How's Runtime PM

# Problem

- **System sleep is not enough to decrease runtime energy consumption**
- **Devices may depend on another device**

# Why's Runtime PM

- **System sleep is not enough to decrease runtime energy consumption**
- **Devices may depend on another device**
- **Solution**
  - Runtime PM

# Runtime PM callback and Helper Functions

● **ksrc/include/linux/pm.h :**

 struct dev_pm_ops {

 …

 int (*runtime_suspend)(struct device *dev);

 int (*runtime_resume)(struct device *dev);

 int (*runtime_idle)(struct device *dev);

 };

● **Ksrc/include/linux/pm_runtime.h**

- pm_runtime_xxx();

# Linux Support for Device Runtime PM

- **Track the number of concurrent users of device**

- **use a per-device reference counter**
  - _get(), _put()

# Runtime PM API: _get(), _put()

- **Tell PM core whether device is in use**

- **I need the device**
  - pm_runtime_get(), _sync(), _noresume()
  - Increment use count, pm_runtime_resume()

- **I'm done**
  - pm_runtime_put(), _sync, _noidle()
  - Decrement use count, pm_runtime_idle()

KESL Kookmin Univ.
Embedded System Lab

# Autosuspend

- **Some device should not allow devices to be suspend**
  - until they have been inactive for some minimum period.
- **A common heuristic**
  - similarly a way of PM qos
- **Init**
  - pm_runtime_set_suspended(dev);
  - pm_runtime_use_autosuspend(dev);
  - pm_runtime_set_autosuspend_delay(dev, DELAY);
- **I need the device**
  - pm_runtime_get_sync(dev);
- **I'm done**
  - pm_runtime_put_sync_autosuspend(dev);

# However, Real world - Runtime PM

- **Exynos SPI driver**
- **Runtime PM often supported after long delay.**

```
1  int s3c64xx_spi_probe(platform_device
       *pdev)
2  {
3    /* allocate controller resources...*/
4    pm_runtime_get_sync(dev);
5    /*initialize the controller...*/
6  }
7
8  int s3c64xx_spi_remove(platform_device
       *pdev)
9  {
10   /* deinitialize controller...*/
11   pm_runtime_put(dev);
12   /* free controller resources...*/
13 }
```

```
1  void s3c64xx_spi_work(work_struct *work)
2  {
3    pm_runtime_get_sync(dev);
4    while (!list_empty(queue)) {
5      /* transmitting message... */
6    }
7    pm_runtime_put(dev);
8  }
9
10 int s3c64xx_spi_setup(spi_device *spi)
11 {
12   pm_runtime_get_sync(dev);
13   /* set up SPI, like tx rate... */
14   pm_runtime_put(dev);
15 }
```

Hand-tuned PM in the Exynos SPI driver

KESL Kookmin Univ. Embedded System Lab

# However, Real world - Runtime PM

- **Complex drivers make it hard to do runtime PM.**
  - It is similar with a Lock.


- **The impact of bad drivers.**
  - Because clock, power, voltage depends on the domains.
  - sharing the same clock source.
  - sharing the same power
  - sharing the same voltage

# Hotplug

# NOHZ_IDLE

KESL Kookmin Univ. Embedded System Lab

# Queue work on power efficient wq

- **Problem**
  - Work-queues can be performance or power oriented.
  - keep them running on a single cpu
    - it remains cache hot
  - but consider big.LITTLE platform

# Queue work on power efficient wq

- **Problem**
  - Work-queues can be performance or power oriented.
  - keep them running on a single cpu
    - it remains cache hot
  - but consider big.LITTLE platform

- **Solution**
  - Power oriented WQ
  - **Give scheduler the liberty**
    - to choose target cpu for running work handler.

- **Power efficient WQ**
  - workqueue is allocated with WQ_UNBOUND flag
  - Idle time for many cpus has increased considerably on **big.LITTLE platform**

# Thermal - iphone, galaxy

# Thermal - G6

# Thermal

**Problem?**

# Thermal

## Spread out

# Today

process/thread -> performance -> multi-core performance scalability-> lock-> bottleneck-> cache cohere system problem-> per-core partition approach

# Next Step.

## Energy-aware scheduling: EAS

1. **CFS scheduler(User level, Tools, Kernel level)**
2. **Load Balancer(Group Scheduling, Bandwidth Control, PELT)**
3. **EAS features**

KESL Kookmin Univ. Embedded System Lab

# 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**

- **Scheduler Activations**
  - https://cgi.cse.unsw.edu.au/~cs3231/12s1/lectures/SchedulerActivations.pdf
- **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**
- **"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**

KESL Kookmin Univ.
Embedded System Lab