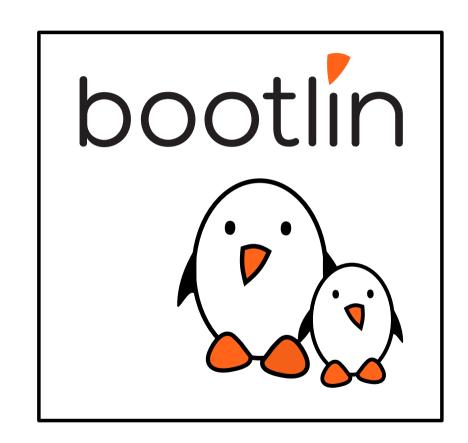


# Linux Power Management Features, Part 2

Théo Lebrun@bootlin.com

Embedded Linux Conference Europe 2025

© Copyright 2004-2025, Bootlin.
Creative Commons BY-SA 3.0 license.
Corrections, suggestions, contributions and translations are welcome!



# P

### Théo Lebrun

- Embedded Linux engineer and trainer at Bootlin
- ▶ Joined Bootlin in 2022, following an internship
  - Linux kernel driver development on embedded systems
    - Suspend-to-RAM for TI J7200 SoC
    - Upstreaming of Mobileye SoCs
  - PipeWire ecosystem
  - Open-source focus
- Current maintainer of elixir.bootlin.com
- Living in Lyon, not-south-nor-north east of France
- theo.lebrun@bootlin.com

# Agenda

Audience: kernel driver developers

Goal: spread the word about existing device PM APIs

### **Covered last year:**

- System-wide suspend
- Runtime power management
- Interactions between the two

### This year:

- 1. How to reason about PM?
- 2. Linux PM features
  - Runtime PM features
  - Power domains
  - PM QoS





- ▶ What is it?
  - Turning off unneeded hardware capabilities
  - Or slowing it down
- ▶ Why we care?
  - Mostly to last longer on batteries
- ► How to do it?
  - It is a tradeoff
  - Either hardware resources are available now
  - Or they are accessible after some latency amount



- Is it complex?
  - Yes
- ► Throttling makes this more than a one dimension axis
  - HW is running at reduced speed
  - Do you keep it as is when a request comes in?
  - Or you pay the transition cost to get full HW capabilities?
- Decisions are platform dependent
  - What can your hardware do?
  - What is most efficient?
  - Slow and steady? Race to halt?
- Decisions are project dependent
  - What is your workload?
  - Can you measure it?
  - What is acceptable?



- Some more issues?
  - Yes
- Making non-platform specific code is hard
  - For a given device, each resource can be in a few different states
  - Platform state is exponentially large
  - Current solution: under-define power states
  - Example: network switch over PCI. What does "it is runtime suspended" mean?
- ► The "we aren't alone" problem
  - Boards aren't single-CPU anymore
  - Sharing resources is hard
  - Don't assume we have full knowledge
  - Don't assume we are making decisions



# Runtime PM



### Runtime PM: introduction

- Per device operations
- Suspend & resume
- Refcount-based
- ► Tree device hierarchy

### include/linux/pm.h

```
// Device API
struct dev_pm_ops {
    /* Device is active but not needed anymore. */
    int (*runtime_suspend)(struct device *dev);
    /* Device is suspended but needed. */
    int (*runtime_resume)(struct device *dev);
    /* ... */
};
```

#### include/linux/pm\_runtime.h

```
// Consumer API
void pm_runtime_enable(struct device *dev);
void pm_runtime_disable(struct device *dev);
int pm_runtime_get(struct device *dev);
int pm_runtime_get_sync(struct device *dev);
int pm_runtime_put(struct device *dev);
int pm_runtime_put(struct device *dev);
int pm_runtime_put_sync(struct device *dev);
```



# Runtime PM: dummy example

```
static int x probe(struct platform device *pdev)
{
        pm runtime enable(&pdev->dev);
        pm runtime get sync(&pdev->dev);
        return 0:
static int x runtime suspend(struct device *dev) { /* ... */ }
static int x_runtime_resume(struct device *dev) { /* ... */ }
static const struct dev pm ops x dev pm ops = {
        RUNTIME PM OPS(x runtime suspend, x runtime resume, NULL)
};
static struct platform driver x driver = {
        .probe = x probe,
        .driver.pm = pm ptr(&x dev pm ops),
};
module_platform_driver(x_driver);
```



### Runtime PM without callbacks

What about runtime PM without callbacks registered?

```
static int x_probe(struct platform_device *pdev)
{
          pm_runtime_enable(&pdev->dev);
          pm_runtime_get_sync(&pdev->dev);
          return 0;
}

static struct platform_driver x_driver = {
               .probe = x_probe,
};
module_platform_driver(x_driver);
```



- First reason: waking up the device model hierarchy
- ► Apart from calling callbacks, runtime PM has other side effects: parent devices might? get resumed as pm\_runtime\_get(dev) calls pm\_runtime\_get(dev->parent)
  - ? not if parent is already active
  - ? not if parent is disabled
  - ? not if parent ignores his children
- ► The parent refcount, dev->power.usage\_count, is always incremented however



- Example: dw\_spi\_pci, a SPI controller driver on the PCI bus.
- ▶ It has no ->runtime suspend|resume() callbacks.

### drivers/spi/spi-dw-pci.c

```
static SIMPLE DEV PM OPS(dw spi pci pm ops, dw spi pci suspend, dw spi pci resume);
static struct pci driver dw spi pci driver = {
        .name = DRIVER NAME,
        .id table = dw spi pci ids,
        .probe = dw spi pci probe,
        .remove = dw spi pci remove,
        .driver
                       = &dw spi pci pm ops,
        },
};
module pci driver(dw spi pci driver);
```



It asks the SPI subsystem to do runtime PM operations on its behalf.

drivers/spi/spi-dw-core.c

```
struct spi controller *host = spi alloc host(dev, 0);
if (!host)
        return - ENOMEM:
// ...
host->max speed hz = dws->max freq;
host->flags = SPI CONTROLLER GPIO SS;
host->auto runtime pm = true;
// ...
ret = spi register controller(host);
if (ret) {
        dev err probe(dev, ret, "problem registering spi host\n");
        goto err dma exit;
```



And finishes its probe by enabling runtime PM.

drivers/spi/spi-dw-pci.c

```
static int dw spi pci probe(struct pci dev *pdev, const struct pci device id *ent)
        // ...
        pm runtime set autosuspend delay(&pdev->dev, 1000);
        pm runtime use autosuspend(&pdev->dev);
        pm runtime put autosuspend(&pdev->dev);
        pm runtime allow(&pdev->dev);
        return 0;
err free irq vectors:
        pci free irq vectors(pdev);
        return ret;
```



#### Conclusion?

- dw\_spi\_pci asks its framework subsystem (SPI) to pm\_runtime\_get|put() automatically on requests.
- This in turn signals to parent devices, ie the PCI bus controller, when the bus can be safely shut down.
- Autosuspend will kick in to delay the suspend. This minimises wasteful suspend/resume cycles during bursts of operations. Parent might have its own autosuspend delay.
- ▶ No ->runtime\_suspend|resume() implementation required inside dw\_spi\_pci.
- Remember the device model is recursive.
  Think GPIO expander over I2C over USB over PCI over platform bus.



# Runtime PM without callbacks: device links

- Second reason: trigger implicit behavior reacting to runtime PM refcount.
- ► That is done through struct device\_link with DL\_FLAG\_PM\_RUNTIME. See also DL\_FLAG\_RPM\_ACTIVE.
- Example usage:

#### drivers/net/phy/phy\_device.c



# Runtime PM without callbacks: device links

- Subsystem using device links with DL\_FLAG\_PM\_RUNTIME are:
  - drivers/net/phy/ for attaching MAC0 to MAC1 if MAC0 uses MAC1's PHY;
  - pinctrl supports linking pin controllers to all their consumers; the flag link\_consumers is used by pinctrl-stmfx.c and stm32/pinctrl-stm32.c;
  - pci for attaching quirked multi-function devices together;
  - dev->dev\_pm\_domain;
  - pmdomain, ie Generic PM Domain;
  - ~36 drivers calling it directly.

I was starting to pull this, and then tried to figure out what the heck "genpd" is.

— Linus



# Power domains



- dev\_pm\_domain is a field inside struct device.
- It got introduced prior to the Git history.
- Of interest to us are dev->dev\_pm\_domain.ops.runtime\_suspend|resume().

### include/linux/pm.h

```
struct dev_pm_domain {
    struct dev_pm_ops ops;
    int (*start)(struct device *dev);
    void (*detach)(struct device *dev, bool power_off);
    int (*activate)(struct device *dev);
    void (*sync)(struct device *dev);
    void (*dismiss)(struct device *dev);
    int (*set_performance_state)(struct device *dev, unsigned int state);
};
```



### Limitations of dev\_pm\_domain:

- One domain per device.
- Not straight-forward to implement.
- Not fitting the object model.
  - The domain is not a device and,
  - No device is marked as providing the domain.
- ▶ It therefore cannot fit in with hardware description (ie devicetree).



- Limitations are addressed by pmdomain, a subsystem that is implemented by piggybacking on dev->dev\_pm\_domain. It used to be called genpd.
- Upstreamed by Rafael J. Wysocki in July 2011.
- Providers must register themselves:

### include/linux/pm\_domain.h



► Consumers aren't expected to use an API (other than runtime suspend/resume); binding is done through devicetree properties & phandles (pinctrl style).

arch/arm64/boot/dts/ti/k3-j7200-mcu-wakeup.dtsi

```
k3 pds: power-controller {
        compatible = "ti,sci-pm-domain";
        #power-domain-cells = <2>;
        bootph-all;
};
wkup vtm0: temperature-sensor@42040000 {
        compatible = "ti,j7200-vtm";
        reg = <0x00 0x42040000 0x00 0x350>,
              <0x00 0x42050000 0x00 0x350>;
        power-domains = <&k3 pds 154 TI SCI PD EXCLUSIVE>;
        #thermal-sensor-cells = <1>;
        bootph-pre-ram;
};
```



```
commit 17f88151ff190b9357f473d7704eee7ae3097d11
Author: Franklin S Cooper Jr <fcooper@ti.com>
```

Date: Mon Sep 11 15:11:44 2017 -0500

i2c: davinci: Add PM Runtime Support

66AK2G has I2C instances that are not apart of the ALWAYS\_ON power domain unlike other Keystone 2 SoCs and OMAPL138. Therefore, pm\_runtime is required to insure the power domain used by the specific I2C instance is properly turned on along with its functional clock.

Signed-off-by: Franklin S Cooper Jr <fcooper@ti.com>

Acked-by: Sekhar Nori <nsekhar@ti.com>

Signed-off-by: Wolfram Sang <wsa@the-dreams.de>



- Mainly, a pmdomain is two callbacks; power-on and power-off.
- ▶ Power-on is called implicitly when going from zero consumer devices active to one.
- Power-off is called whenever all devices inside the domain are runtime suspended.

#### include/linux/pm\_domain.h

Benefits of pmdomain compared to raw dev->dev\_pm\_domain:

- Straight forward to implement (example: imx/imx93-pd.c).
- Fitting the object model: provider is an identified device.
- ► Mark resources inside devicetree (as it should be).
- Still one domain per device but we can implement domain hierarchy.
  - Kevin Hilman has a series to add support for it in DT.
  - Issue with flat firmware-provided PM domains, as scmi.
- Zero code inside ->runtime\_suspend|resume() for consumers.
  - Transparent to consumer drivers, if they do runtime PM.



### Some more pmdomain features:

- ► GENPD\_FLAG\_PM\_CLK / pm\_clk\_\* infrastructure. Attaches clocks to a given power domain, and those are enabled/prepared implicitly.
- ► Callbacks when devices get attached. Example usage: drivers use the PM\_CLK infrastructure to attach device clocks coming from devicetree.

```
int (*attach_dev)(struct generic_pm_domain *domain, struct device *dev);
void (*detach_dev)(struct generic_pm_domain *domain, struct device *dev);
```

- ► GENPD\_FLAG\_ACTIVE\_WAKEUP for handling power domains on the system-wide suspend wakeup path.
- ► GENPD\_FLAG\_RPM\_ALWAYS\_ON for some (broken?) platforms.



# PM QoS



### PM QoS: introduction

- ▶ PM QoS (*Quality of Service*) is about registering performance expectations.
- ▶ Per-device PM QoS is about **requests** to respect a maximum resume latency.
- ► All requests are aggregated into a single resume latency to respect.

#### include/linux/pm\_qos.h

```
enum dev pm gos reg type {
       DEV PM QOS RESUME LATENCY = 1, // set max PM resume latency
       DEV PM QOS LATENCY TOLERANCE, // value interpreted by driver
       DEV PM QOS MIN FREQUENCY, // passed to cpufreq
       DEV PM QOS MAX FREQUENCY, // passed to cpufreq
       DEV PM QOS FLAGS,
};
int dev pm gos add request(struct device *dev, struct dev pm gos request *req,
                          enum dev pm qos req type type, s32 value);
int dev pm qos update request(struct dev pm qos request *req, s32 new value);
int dev_pm_qos_remove_request(struct dev_pm_qos_request *req);
s32 dev pm gos read value(struct device *dev, enum dev pm gos reg type type);
```

# PM QoS: impact on runtime suspend

- If resume\_latency is zero, never runtime suspend; rpm\_check\_suspend\_allowed().
- ▶ Otherwise, we might? runtime suspend if we know the operation is quick enough.
  - That is coming from past experience: how slow were the worst runtime suspend/resume?
  - ? only if device belongs to PM domain
  - ? only if PM domain has a governor
- See default\_suspend\_ok() in drivers/pmdomain/governor.c.

```
dev.effective_constraint_ns = dev.resume_latency
for child_dev in dev.children:
    if child_dev.effective_constraint_ns < dev.effective_constraint_ns:
        effective_constraint_ns = child_dev.effective_constraint_ns

# worst past measurements
dev.effective_constraint_ns -= dev.suspend_latency_ns
dev.effective_constraint_ns -= dev.resume_latency_ns

if dev.effective_constraint_ns >= 0:
    dev.suspend()
```



# PM QoS: impact on runtime suspend

- suspend\_latency\_ns is time for dev->runtime\_suspend() and pmdomain->stop(dev).
- resume\_latency\_ns is time for pmdomain->stop(dev) and dev->runtime\_suspend().
- ► **Notice:** this does not include the pmdomain->power\_on() and pmdomain->power\_off() duration, as it isn't per-device.
  - Depending on the platform, it might be the biggest time sink.



# PM QoS: PM domains

- ► PM domains QoS is done through separate measurements: power\_on\_latency\_ns and power\_off\_latency\_ns.
  - Behavior is similar to suspend\_latency\_ns and resume\_latency\_ns.
  - Accounts for PM subdomains.
  - Accounts for children devices known suspend/resume timings.
- ► Values can come from devicetree; in that case, further ->power\_on|off() are not timed.
- See \_\_default\_power\_down\_ok() in drivers/pmdomain/governor.c.



# Conclusion



# More features (in bulk)

- cpufreq configures CPU frequencies; it gets configured through QoS APIs.
- cpuidle picks from its list of idle states the deepest acceptable state.
  See Monday talk by Dhruva Gole (TI) and Kevin Hilman (BayLibre).
- ▶ PM domains provide OPP: "The set of discrete tuples consisting of frequency and voltage pairs that the device will support per domain are called Operating Performance Points". Documentation. Contributed in 2010 by Nishanth Menon (TI).
- ► PM domains have set\_hwmode\_dev() / get\_hwmode\_dev() for toggling hardware controlled PM. Only one user (Qualcomm GDSC).
- ▶ Related ongoing work: The Case for an SoC Power Management Driver Stephen Boyd, Google (2024).

# Thank you!

Questions?

Théo Lebrun theo.lebrun@bootlin.com

Slides under CC-BY-SA 3.0

https://bootlin.com/pub/conferences/