# Power Management Using PM Domains on SH7372

Rafael J. Wysocki

Renesas / Faculty of Physics U. Warsaw / SUSE Labs

October 3, 2011

1 / 18

### SH7372 Characteristics

- System on a Chip (SoC).
- Two independent CPU cores (ARM and SuperH).
- Possibility to manipulate device clocks (enable/disable).
- Power domains that can be turned on and off through register writes.
- Hierarchy of power domains.
- No BIOS (description through board files in the kernel).

## SH7372 Characteristics

- System on a Chip (SoC).
- Two independent CPU cores (ARM and SuperH).
- Possibility to manipulate device clocks (enable/disable).
- Power domains that can be turned on and off through register writes.
- Hierarchy of power domains.
- No BIOS (description through board files in the kernel).

#### Simple power management model

- Device on (full-power state).
- Clocks off (low-power state, level 1).
- Power removed from power domain (low-power state, level 2).
- Wakeup latency constraints determine the state to choose (run time).



# Power Domains On SH7372

```
C5: base (mother) domain, CPG, KEYSC, CMT, RWDT, GPIO
```

A4LC: LCDC, DSI, MERAM (video)

A4MP: SPU2, FSI (audio)

D4: ARM debug

A4R: SH4AL-DSP, INTCS, DMAC, IIC, TMU, MSIOF, CMT0,

CEU, CSI (SH CPU core, I/O)

A3RI: ISP (camera capture unit)

A3RV: VPU (video encode/decode unit)

A4S: INTCA, MFI, SBSC (interrupt and SDRAM controllers)

A3SG: SGX (3D graphics)

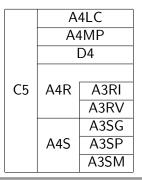
A3SP: SCIF, MSIOF, IIC, USB, SDHI, MMCIF, HDMI (I/O)

A3SM: ARM Cortex-A8 CPU core

# Hierarchy Of Power Domains

	A4LC A4MP D4	
C5		
	A4R	A3RI
		A3RV
	A4S	A3SG
		A3SP
		A3SM

# Hierarchy Of Power Domains



It turns out that A3RV depends on A4LC (because of MERAM).

# Consequences Of The Design

- Every device is a direct member of one power domain.
- One power domain may have multiple masters (e.g. A3RV).
- It is desirable to turn off A4R when A3RI and A3RV are off.
- It is desirable to turn off A4S when A3SG, A3SP, A3SM are off.

# Consequences Of The Design

- Every device is a direct member of one power domain.
- One power domain may have multiple masters (e.g. A3RV).
- It is desirable to turn off A4R when A3RI and A3RV are off.
- It is desirable to turn off A4S when A3SG, A3SP, A3SM are off.

#### However

There are devices with cross-domain dependencies (e.g. INTCA, DMAC).

# Consequences Of The Design

- Every device is a direct member of one power domain.
- One power domain may have multiple masters (e.g. A3RV).
- It is desirable to turn off A4R when A3RI and A3RV are off.
- It is desirable to turn off A4S when A3SG, A3SP, A3SM are off.

#### However

There are devices with cross-domain dependencies (e.g. INTCA, DMAC).

The I/O runtime PM framework had to be extended so that the dependencies related to power domains could be handled.

### **CPUidle**

### Theory: Put idle CPUs into low-power states (no code execution)

- CPU scheduler knows when a CPU is idle.
- Next usage time information from clock events.
- Maximum acceptable wakeup latency from PM QoS.
- CPU low-power states (C-states) characteristics are known.
- Governors decide what state to go for.

### **CPUidle**

### Theory: Put idle CPUs into low-power states (no code execution)

- CPU scheduler knows when a CPU is idle.
- Next usage time information from clock events.
- Maximum acceptable wakeup latency from PM QoS.
- CPU low-power states (C-states) characteristics are known.
- Governors decide what state to go for.

#### Practice: Not so easy

- Power domain A4S is shared with I/O devices.
- Deep "C-states" may involve removing power from A4S.
- How much time is it going to take to restore power (wakeup)?
- What is the power break even time?



## Runtime PM Framework

### Turning devices off (when idle) and on

- The core:
  - Handles concurrency (locking etc.).
  - Takes care of device dependencies (parents vs children).
  - Provides reference counting facilities (detection of idleness).
  - Provides common helpers (e.g. pm\_runtime\_suspend()).
- Subsystems and drivers:
  - Provide callbacks.
  - 4 Handle wakeup events (remote wakeup).

## Runtime PM Framework

### Turning devices off (when idle) and on

- The core:
  - Handles concurrency (locking etc.).
  - Takes care of device dependencies (parents vs children).
  - Provides reference counting facilities (detection of idleness).
  - Provides common helpers (e.g. pm\_runtime\_suspend()).
- Subsystems and drivers:
  - Provide callbacks.
  - 4 Handle wakeup events (remote wakeup).

Power domains have to be taken into account.

## Runtime PM Framework

#### Turning devices off (when idle) and on

- The core:
  - Handles concurrency (locking etc.).
  - 2 Takes care of device dependencies (parents vs children).
  - Provides reference counting facilities (detection of idleness).
  - Provides common helpers (e.g. pm\_runtime\_suspend()).
- Subsystems and drivers:
  - Provide callbacks.
  - 4 Handle wakeup events (remote wakeup).

Power domains have to be taken into account.

Subsystem callbacks may be overriden by power management domain callbacks (representation via struct dev\_pm\_domain).

## Generic PM Domains

### Simple framework for representing power domains

- Stop/start operations for member devices (may correspond to enabling/disabling clocks).
- Save state/restore state operations for member devices (represented by drivers' .runtime\_suspend() and .runtime\_resume() callbacks).
- Power off/power on operations for entire domains.
- Domain hierarchies.

### Generic PM Domains

### Simple framework for representing power domains

- Stop/start operations for member devices (may correspond to enabling/disabling clocks).
- Save state/restore state operations for member devices (represented by drivers' .runtime\_suspend() and .runtime\_resume() callbacks).
- Power off/power on operations for entire domains.
- Domain hierarchies.

Representation via struct generic\_pm\_domain (new in the 3.1-rc kernels, improvements on their way into the 3.2 kernel).

# Device Suspend

- Device usage counter is 0, the PM core runs rpm\_suspend(dev).
- pm\_genpd\_runtime\_suspend(dev) is called.
- **1** The "stop" operation is carried out for the device.
- pm\_genpd\_poweroff(pd) is called for the device's domain.
- If all devices in the domain are stopped and its subdomains are "off":
  - The states of all devices in the domain are saved.
  - 2 The "power off" operation is carried out for the domain.

# Device Suspend

- Oevice usage counter is 0, the PM core runs rpm\_suspend(dev).
- pm\_genpd\_runtime\_suspend(dev) is called.
- The "stop" operation is carried out for the device.
- pm\_genpd\_poweroff(pd) is called for the device's domain.
- If all devices in the domain are stopped and its subdomains are "off":
  - The states of all devices in the domain are saved.
  - 2 The "power off" operation is carried out for the domain.

There are multiple conditions that have to be satisfied for this procedure to be completed (it may be aborted at any point).

# Device Suspend

- Oevice usage counter is 0, the PM core runs rpm\_suspend(dev).
- pm\_genpd\_runtime\_suspend(dev) is called.
- The "stop" operation is carried out for the device.
- pm\_genpd\_poweroff(pd) is called for the device's domain.
- If all devices in the domain are stopped and its subdomains are "off":
  - The states of all devices in the domain are saved.
  - 2 The "power off" operation is carried out for the domain.

There are multiple conditions that have to be satisfied for this procedure to be completed (it may be aborted at any point).

Timing (PM QoS) constraints should be taken into account when deciding whether or not to carry out "power off" (not implemented yet).

### Device Resume

- Device usage counter is incremented, the PM core runs rpm\_resume(dev).
- pm\_genpd\_runtime\_resume(dev) is called.
- If necessary, the "power on" operation is carried out for the device's domain.
  - This has to abort all instances of pm\_genpd\_poweroff(pd) running for the same domain.
  - It has to be done recursively for all of the "master" domains before.
- If necessary, the device's state is restored.
- The "start" operation is carried out for the device.

### Device Resume

- Device usage counter is incremented, the PM core runs rpm\_resume(dev).
- pm\_genpd\_runtime\_resume(dev) is called.
- If necessary, the "power on" operation is carried out for the device's domain.
  - This has to abort all instances of pm\_genpd\_poweroff(pd) running for the same domain.
  - It has to be done recursively for all of the "master" domains before.
- If necessary, the device's state is restored.
- The "start" operation is carried out for the device.

There are devices with cross-domain dependencies requiring special handling (interrupt controller, DMA engine). Not supported yet.

# Power Management Quality Of Service (PM QoS)

#### Observations

- When a device is suspended, the kernel has to decide whether or not to turn its power domain (or its master) off.
- It takes time to turn power domains off/on and to restore devices' registers.
- There may be requirements regarding the time devices can spend in the non-working state.
- The device suspend code should take those requirements into account.

# Power Management Quality Of Service (PM QoS)

#### Observations

- When a device is suspended, the kernel has to decide whether or not to turn its power domain (or its master) off.
- It takes time to turn power domains off/on and to restore devices' registers.
- There may be requirements regarding the time devices can spend in the non-working state.
- The device suspend code should take those requirements into account.

#### Device PM QoS

Framework allowing kernel subsystems to specify wakeup latency constraints for I/O devices (scheduled for inclusion into the 3.2 kernel).

# Power Management Quality Of Service (PM QoS)

#### Observations

- When a device is suspended, the kernel has to decide whether or not to turn its power domain (or its master) off.
- It takes time to turn power domains off/on and to restore devices' registers.
- There may be requirements regarding the time devices can spend in the non-working state.
- The device suspend code should take those requirements into account.

#### Device PM QoS

Framework allowing kernel subsystems to specify wakeup latency constraints for I/O devices (scheduled for inclusion into the 3.2 kernel).

The generic PM domains code will take PM QoS constraints into account.

Well, no ACPI BIOS, no system sleep states, right?

Well, no ACPI BIOS, no system sleep states, right?

### Life's more complicated than that

- Android uses system suspend for power saving.
- Wakeup from deep low-power states may require special signaling.

Well, no ACPI BIOS, no system sleep states, right?

### Life's more complicated than that

- Android uses system suspend for power saving.
- Wakeup from deep low-power states may require special signaling.

Difference between "sleep" and "power off"

The preservation of RAM contents.

Well, no ACPI BIOS, no system sleep states, right?

#### Life's more complicated than that

- 4 Android uses system suspend for power saving.
- Wakeup from deep low-power states may require special signaling.

#### Difference between "sleep" and "power off"

The preservation of RAM contents.

#### Still

SH7372 doesn't provide any "enter a sleep state" logic.

# "Sleep" Vs "Low-Power"

Sleep states are special low-power states of the whole system.

# "Sleep" Vs "Low-Power"

Sleep states are special low-power states of the whole system.

#### "Sleep" state

- Reduced energy consumption.
- User space is not being run.
- Limited set of devices whose wakeup signals cause user space to be "thawed" (controlled by user space).
- System may be instructed to go into it at any time.

# "Sleep" Vs "Low-Power"

Sleep states are special low-power states of the whole system.

#### "Sleep" state

- Reduced energy consumption.
- User space is not being run.
- Limited set of devices whose wakeup signals cause user space to be "thawed" (controlled by user space).
- System may be instructed to go into it at any time.

There need to be special procedures leading to and from a sleep state

• System suspend and resume, respectively.

# System Suspend And Power Domains

#### Rules

- The state of devices in a power domain has to be saved before removing power from it.
- Oomains that were "off" before system suspend generally should stay this way.

# System Suspend And Power Domains

#### Rules

- The state of devices in a power domain has to be saved before removing power from it.
- ② Domains that were "off" before system suspend generally should stay this way.

#### Wakeup device problem

- A power domain in the "off" state contains a device whose input lines can generate wakeup events.
- ② User space doesn't want to be woken up by this device.

# System Suspend And Power Domains

#### Rules

- The state of devices in a power domain has to be saved before removing power from it.
- ② Domains that were "off" before system suspend generally should stay this way.

#### Wakeup device problem

- A power domain in the "off" state contains a device whose input lines can generate wakeup events.
- User space doesn't want to be woken up by this device.

In that case it is necessary to restore power to the domain and to reprogram the device during system suspend.



- SoCs (like SH7372) add complexity to power management (power domains, direct control of device clocks).
- @ Generic PM code has to be extended to address the added complexity.
- That code may be modeled on the basis of the SH7372 design.
- The fact that the new code can be tested on real hardware matching basic assumptions closely helps a lot.
- The introduction of new generic PM code allows other platforms to use it without adding code duplication and "reinventing the wheel".
- Discussions with developers working on other platforms resulted in better code.

## References

- R. J. Wysocki, Runtime PM vs System Sleep (http://www.linuxplumbersconf.org/2011/ocw/system/presentations/27/original/system\_sleep\_vs\_runtime\_PM.pdf).
- R. J. Wysocki, Runtime Power Management vs System Sleep (http://events.linuxfoundation.org/slides/2011/linuxcon-japan/lcj2011\_wysocki.pdf).
- R. J. Wysocki, Runtime Power Management Framework for I/O Devices in the Linux Kernel (http://events.linuxfoundation.org/slides/2010/linuxcon2010\_wysocki.pdf).
- R. J. Wysocki, Runtime Power Management Framework for I/O Devices in the Linux Kernel (http://events.linuxfoundation.org/slides/2011/linuxcon-japan/lcj2011\_wysocki2.pdf).

## Documentation And Source Code

- Documentation/power/devices.txt
- Documentation/power/runtime\_pm.txt
- include/linux/cpuidle.h
- include/linux/device.h
- include/linux/pm.h
- include/linux/pm\_domain.h
- include/linux/pm\_runtime.h
- include/linux/pm\_wakeup.h
- include/linux/suspend.h
- o drivers/base/power/
- drivers/cpuidle/
- drivers/cpufreq/
- kernel/power/



# Thanks!

Thank you for attention!



## Thanks!

#### Thank you for attention!

Special thanks to

Renesas Electronics Corp.

Faculty of Physics U. Warsaw

