

PM Infrastructure in the Linux* Kernel

Current Status And Future

Rafael J. Wysocki

Intel Open Source Technology Center

October 4, 2016

Outline

- 1 Introduction
 - The Goal
 - Power Management Frameworks
- 2 System-Wide Power Management (Sleep States)
 - How It Works
 - The Future
- 3 Working-State Power Management
 - I/O Device Runtime PM
 - CPU Power Management
- 4 Resources

What We Are After

Objective

Use only as much energy as needed to achieve sufficient performance.

What software can do

- 1 Determine how much throughput/capacity is needed.
- 2 Determine how much latency is acceptable.
- 3 Enumerate the hardware PM features.
- 4 Use them to deliver exactly as much as necessary.

What about the Linux kernel?

The Kernel's Role

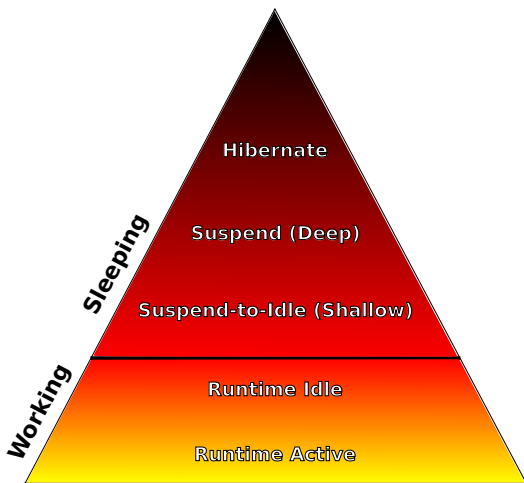
Provide means of control

- State selection (component level).
- Carrying out transitions between states (the “mechanics”).
- Response to events (e.g. wakeup).

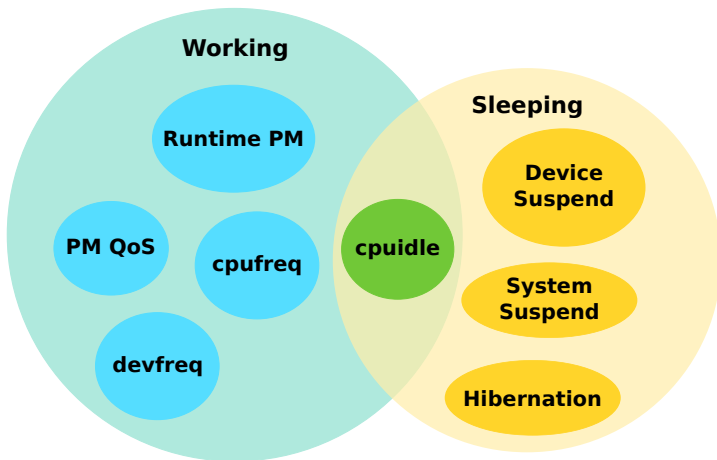
Estimate the needs

- Use information available internally (e.g. from the CPU scheduler).
- React to the actions of user space.
- Follow trends/patterns.

Working-State PM and System-Wide PM (Sleep States)



Power Management Frameworks in The Linux Kernel



System-Wide PM Overview

Global energy-saving states, “frozen” user space

Suspend-to-Idle, Standby, Suspend-to-RAM, Hibernate.

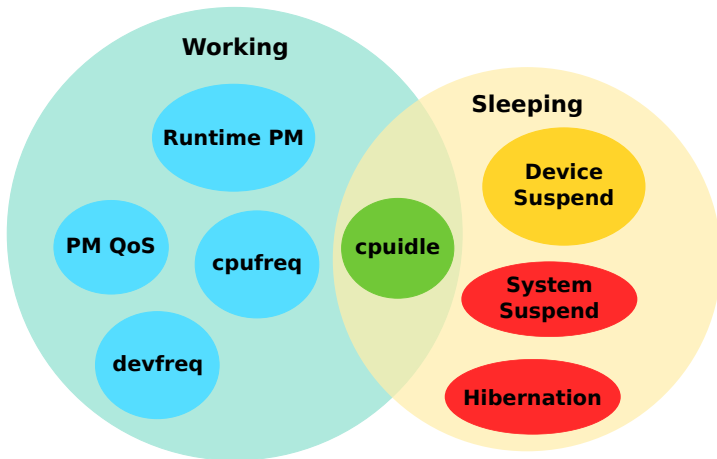
Controlled by user space

- 1 User space selects the target state.
- 2 User space decides when to start transitions.
 - Direct command (*sysfs* write).
 - Autosleep interface.

Used on many systems

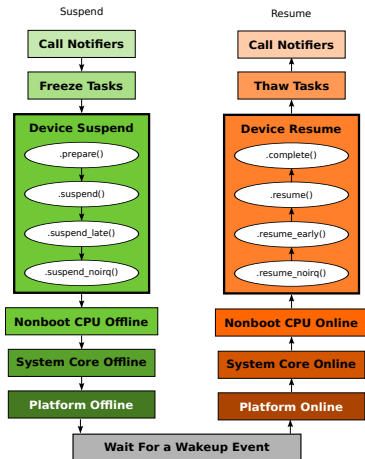
Desktop distributions, Android (autosleep interface).

System-Wide PM Frameworks

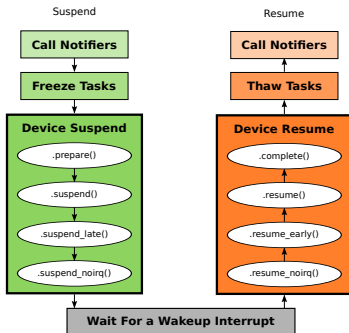


System Suspend/Resume Control Flow

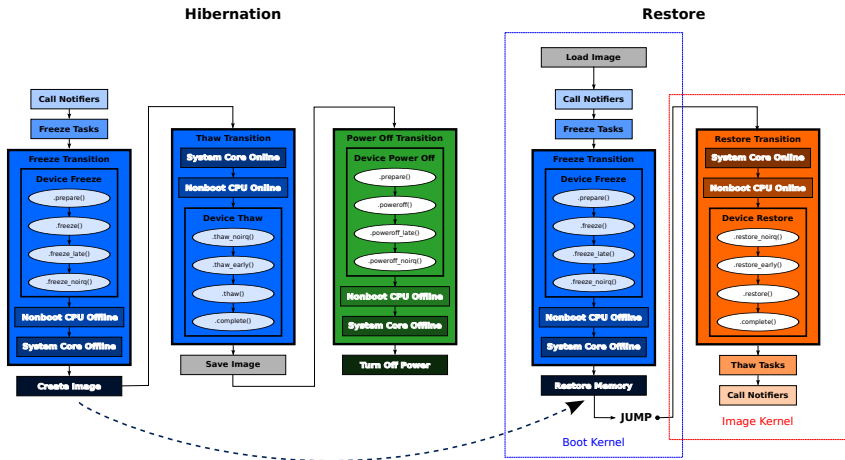
Full Suspend



Suspend to Idle



System Hibernation/Restore Control Flow



The Future of System-Wide PM

Hibernation

- Doesn't go away (supported on ARM64 now, works with KASLR).
- Encrypted images problematic.
- Will persistent memory make it obsolete?

Future platforms may not support “platform offline”

- Suspend-to-RAM and Standby may not make sense.
- CPU offline/core offline steps unnecessary.
- Suspend-to-Idle can leverage the existing infrastructure.

Suspend-to-Idle vs Working-State PM

For long inactivity periods Suspend-to-Idle **always** uses less energy

- Suspends timekeeping (no timer interrupts).
- Longer average time between wakeups.

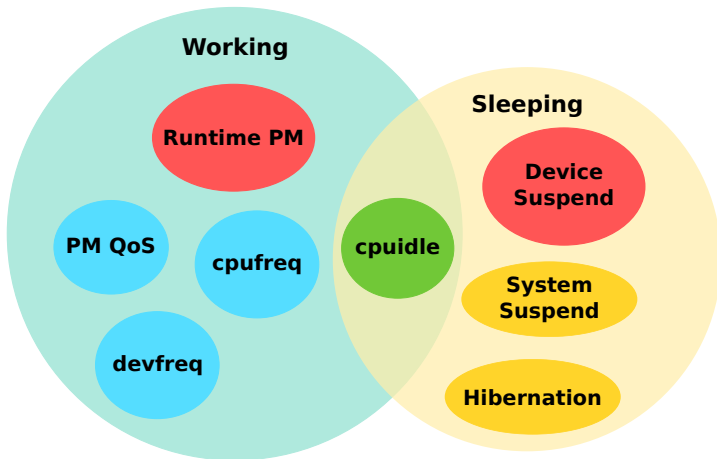
Different mechanisms for different use cases

- Suspend on closing laptop lid (working-state PM insufficient).
- Opportunistic idle (Suspend-to-Idle insufficient).

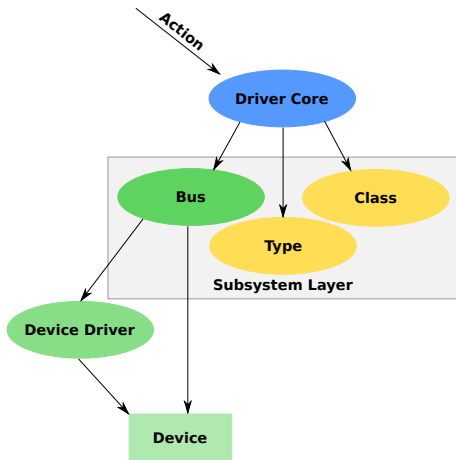
Suspend-to-Idle might be more efficient than it is today

Framework/driver issues to fix.

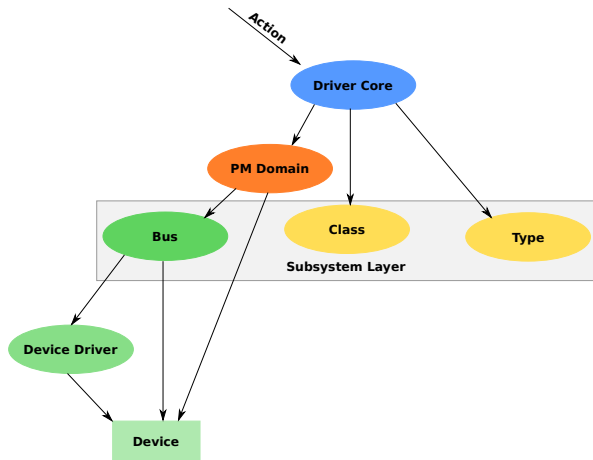
Suspend/Resume of I/O Devices



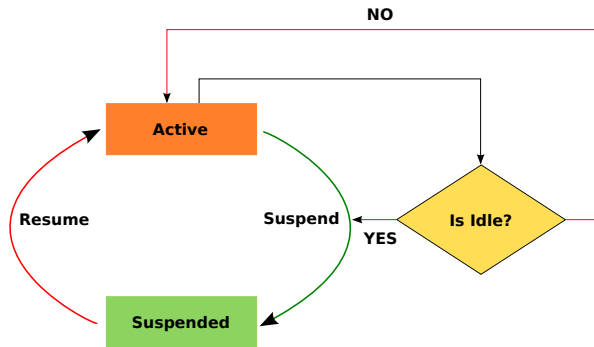
Device PM Operations and The Driver Core



Device PM Operations and PM Domains



Device Runtime PM Operations



Runtime PM on Future Systems

Hardware design trends

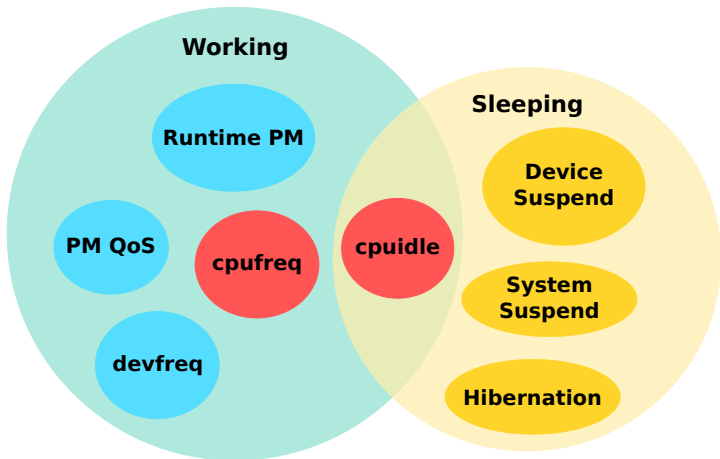
- Increasing integration of components.
- Increasing level of support for aggressive PM.

System-on-a-Chip (SoC) configurations

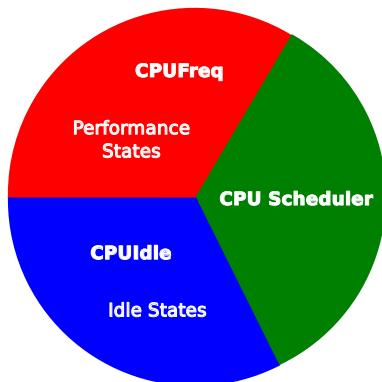
- CPU packages contain I/O devices.
- Package low-power states depend on I/O devices.
- PM features of different components are interdependent.

Challenge: Take dependencies into account

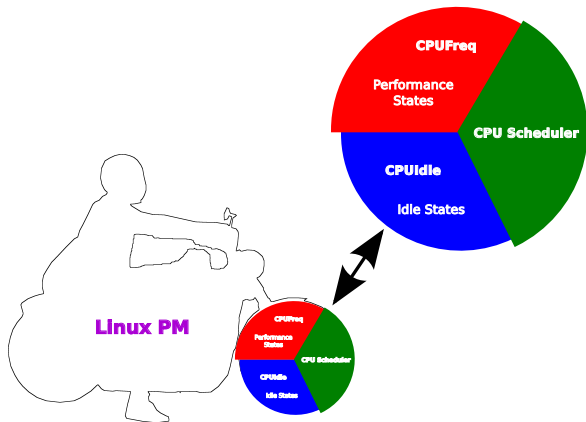
CPU Power Management Frameworks



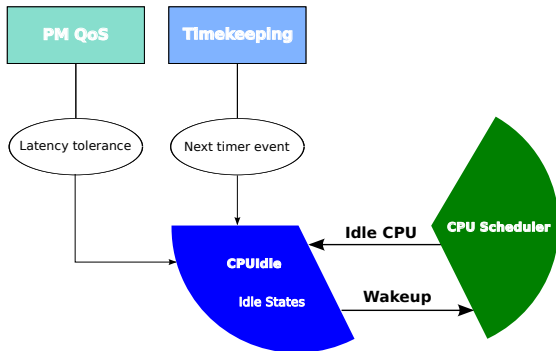
CPU Power Management Overview



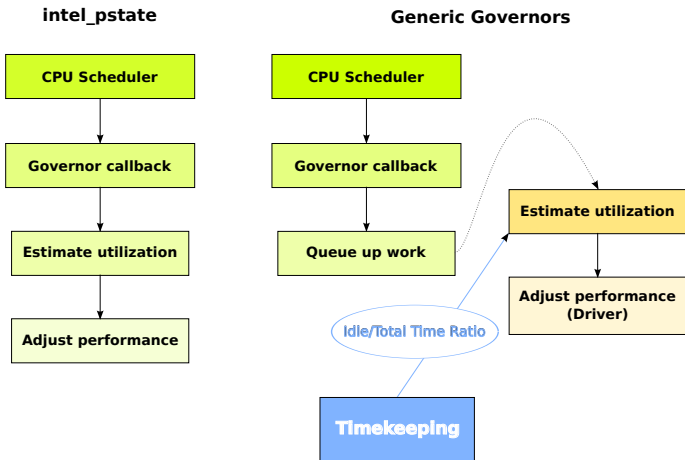
Cooperation Between Components Is Key



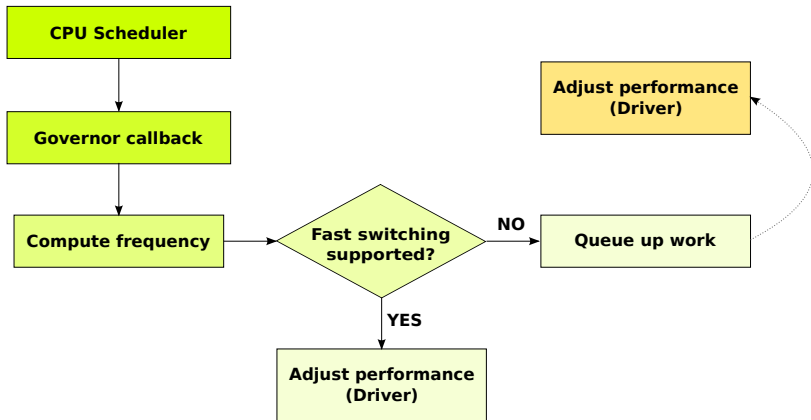
CPUIidle Workflow



CPUFreq: Old-Style Governors and intel_pstate



CPUFreq: schedutil Governor



CPUFreq: Scheduler Hints, Cross-CPU Updates

Observations

- It is good to take blocking on I/O (*IOwait*) into account.
- Different scheduling classes require different handling.

Solution: Scheduler hints (work in progress)

- Pass hints from the CPU scheduler to governor callbacks.
- A hint can represent the reason update or similar.

Observation

In some cases it is beneficial to invoke governor callbacks cross-CPU.

The Future of CPU Power Management

Complex topologies (hardware threads, modules, packages)

Many different scheduling strategies are potentially viable.

PM-Aware Scheduling Conjecture

Energy efficiency may be improved **without hurting performance** by making the CPU scheduler drive CPU power management as a whole.

Predictions (usual disclaimers apply)

- PM-aware scheduling will enter the mainline kernel.
- CPUFreq and CPUIidle will be combined.

Conclusion

- The Linux kernel supports power management in a number of ways.
- Both system-wide and working state (runtime) PM are supported.
- Support for system-wide PM in device drivers is generally better.
- I/O device runtime PM support is improving.
- CPU PM is well supported, consolidation in progress.
- Hardware design trends increase PM complexity.
- Interdependencies between PM features are challenging.

Questions?

References



R. J. Wysocki, *CPUfreq and The Scheduler: Revolution in CPU Power Management* (http://events.linuxfoundation.org/sites/events/files/slides/cpufreq_and_scheduler_0.pdf).



Neil Brown, *Improvements in CPU frequency management* (<http://lwn.net/Articles/682391/>).



Len Brown, *Suspend/Resume at the Speed of Light* (<http://events.linuxfoundation.org/sites/events/files/slides/Brown-Linux-Suspend-at-Speed-of-Light-LC-EU-2015.pdf>).



R. J. Wysocki, *Getting More Out Of System Suspend In Linux* (http://events.linuxfoundation.org/sites/events/files/slides/linux_suspend.pdf).



R. J. Wysocki, *Power Management in the Linux Kernel – Current Status and Future* (http://events.linuxfoundation.org/sites/events/files/slides/kernel_PM_plain.pdf).



Jonathan Corbet, *The cpuidle subsystem* (<http://lwn.net/Articles/384146/>).



R. J. Wysocki, *Why We Need More Device Power Management Callbacks* (https://events.linuxfoundation.org/images/stories/pdf/lfcs2012_wysocki.pdf).

Documentation and Source Code

- Documentation/cpu-freq/*
- Documentation/cpuidle/*
- Documentation/power/devices.txt
- Documentation/power/pci.txt
- Documentation/power/runtime_pm.txt
- include/linux/cpufreq.h
- include/linux/cpuidle.h
- include/linux/device.h
- include/linux/pm.h
- include/linux/pm_runtime.h
- include/linux/suspend.h
- drivers/acpi/processor_idle.c
- drivers/base/power/*
- drivers/cpufreq/*
- drivers/cpuidle/*
- kernel/power/*

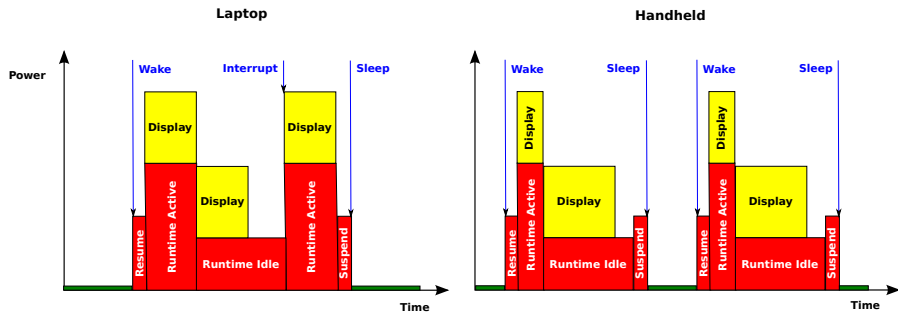
Legal Information

Intel is a trademark of Intel Corporation in the U. S. and other countries.
*Other names and brands may be claimed as the property of others.
Copyright © 2016 Intel Corporation, All rights reserved.

Thanks!

Thank you for attention!

Laptop/Handheld Usage Scenarios



Dark Resume Scenario

