# Linux Kernel Power Management (PM) Framework for ARM 64-bit Processors

L.Pieralisi

21/8/2014 - LinuxCon North America 2014

**ARM**

# Outline

- ARM 32-bit Linux kernel power management support for huge legacy of ARM processors
  - from v4 Uniprocessor kernels to ARM v7 SMP multicluster systems
- Lack of established firmware interfaces is preventing merge of power management software in the mainline kernel
  - Lots of tricky platform specific code, maintained as out of tree BSP branches
  - Power management HW dependent SW layers incompatible with most of kernel SW interfaces
- ARM64 kernel port represents an opportunity to start afresh and learn lessons from the past
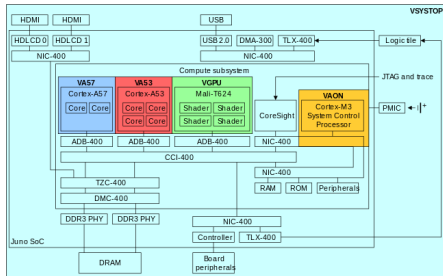
**ARM**

# Legacy ARM 32-bit Linux Kernel Power Management

- Standards ? No, thanks
- ARM vendors power controllers
  - Differentiating through design quirks
- Standard power down procedures started appearing on ARM v7 TRMs
  - What works on ARM v7, might not on v6, and viceversa
  - ARM v7 power down procedure is standard, except for when it is not
  - ....and then came SMP systems....
- Kernel developers left to their own devices
  - Reverse engineering (lack of publicly available specs)
  - Power management interfaces designed for UP systems did not work on SMP
  - Secure/non-secure split overriden in most upstream design
  - Kernel subsystems (eg CPUidle) redesigned to support HW bugs (eg couple idle states)

**ARM**

# ARM Chips Power Management Configurations

- Multiple power islands
  - Devices and CPUs power domains
- Core gating and Cluster gating
  - RAM retention states
  - Caches management
  - Coherency management
  - Locking
- Graphics power domain
  - Integration with CPU states
- System sleep states
- Always-on power domain



**ARM**

# ARM Power Management HW/SW stack (1/2)

- Cross-divisional collaborative effort between ARM SW/HW engineers
- HW prototyping effort (TC2/Juno)
  - Power controller specifications
  - SW stack (Linux kernel/Trusted Firmware/power controller firmware)
- Tackle ARM vendors power control fragmentation
  - Power management HW still a prominent feature of vendors design
  - Design quirks are the rule, not the exception
- Ongoing effort at ARM to standardize main HW control feaures
  - HW SMP coherency management
  - Cache clean and invalidate (SW/HW)
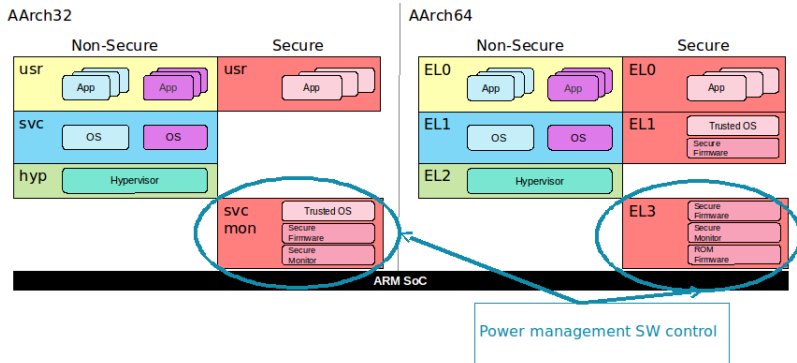  - Wake-up capabilities and GIC design

**ARM**

# ARM Power Management HW/SW stack (2/2)

- Development and deployment of PSCI (Power State Coordination Interface) firmware interface
- Device tree and ACPI bindings standardization
    - Provide OS with standard interfaces
    - Configuration data for power management subsystems (CPU/device idle states)
    - Run-time services (CPU/devices power management)

**Our Goal**

Provide a HW/SW environment to foster power management standardization

**ARM**

# ARM Power Management: Secure Operations



- Power management SW has to deal with secure operations
  - Coherency management, power control commands
  - Need for a standard API to communicate to secure world

# ARM Power Management: Secure Operations

"... I guess what this ultimately comes down to is that we _did_ accept the work-arounds into the kernel source for secure parts - had we not done that and set a clear message like ARM64 does that work-arounds must be done prior to calling the kernel (irrespective of how that happens, iow, whether boot or resume) then we wouldn't have this problem right now. Such a statement would have raised lots of objections though, but with hindsight, it would have been the right thing to do to overrule those objections and just mandate it. It would've been a pain for some people, but we would not be in this situation now where there is no proper solution which works for everyone."

RMK

http://lists.infradead.org/pipermail/linux-arm-kernel/2014-July/268718.html
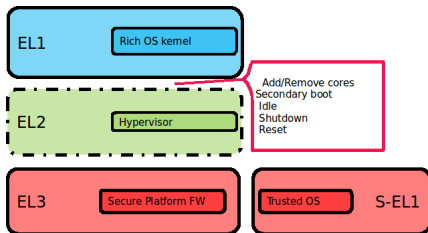
ARM

# ARM Power Management: Early Devices

```
1   int __init ve_spc_init(void __iomem *baseaddr,
2                          u32 a15_clusid, int irq)
3   {
4           int ret;
5           info = kzalloc(sizeof(*info), GFP_KERNEL);
6           if (!info) {
7                   return -ENOMEM;
8           }
9
10          info->baseaddr = baseaddr;
11          info->a15_clusid = a15_clusid;
12
13          [...]
14
15          sync_cache_w(info);
16          sync_cache_w(&info);
17
18          return 0;
19  }
```

```
1   void ve_spc_set_resume_addr(u32 cluster, u32 cpu,
2                               u32 addr)
3   {
4           void __iomem *baseaddr;
5
6           if (cluster >= MAX_CLUSTERS)
7                   return;
8
9           if (cluster_is_a15(cluster))
10                  baseaddr = info->baseaddr +
11                             A15_BX_ADDR0 + (cpu << 2);
12          else
13                  baseaddr = info->baseaddr +
14                             A7_BX_ADDR0 + (cpu << 2);
15
16          writel_relaxed(addr, baseaddr);
17  }
```

- Need to probe device drivers early (before `early_initcalls`)
    - Power controllers drivers
    - CPU memory mapped peripherals (no access through coprocessor - need DT probing)
- Resulting code incompatible with kernel driver model
    - Code exiled to arch/arm, because it does not belong anywhere else

ARM

# Power State Coordination Interface



- Defines a standard interface for making power management requests across exception levels/operating systems
- Support virtualisation and a communication channel between normal and secure world
- Allow secure firmware to arbitrate power management requests from secure and non-secure software
- Default method for power control in Linux ARM64 kernel
    - CPU on/off (hotplug, cold-boot)
    - CPU suspend (suspend to RAM, CPUidle)
- Spec available at
  http://infocenter.arm.com/help/index.jsp?topic=/com.arm.doc.
  den0022b/index.html

# ARM64 Linux Power Management Requirements

- Reliance on PSCI as firmware interface
- Trusted firmware PSCI implementation publicly available for FVP models and Juno
  `https://github.com/ARM-software/arm-trusted-firmware`
- Device tree and ACPI bindings to standardize control data and improve code generality
  - Take part in ACPI standardization effort `http://www.acpi.info`
  - Actively push for Device Tree bindings for power management SW components
- Push back on design quirks and related SW hacks
- More collaboration between ARM and ARM partners on HW/SW interfaces
  - HW standardization effort must be deployed in ARM vendors designs
  - SW interfaces adopted by ARM partners Linux kernel developers

**ARM**

# Linux Kernel Power Management Subsystems

- CPUidle: framework managing idle threads
  - CPUs enter idle states when idling, with power depth levels that depend on power management HW
- CPU hotplug: removing a CPU from a running system
  - Not designed for power management, abused as such
- Runtime PM
  - Devices power management
- Suspend to RAM: system wide (CPU and devices) sleep states
  - Triggered from userspace, also used as autosuspend when system is idle (eg on Android)
- Suspend to disk (aka Hibernation): same as Suspend-to-RAM, system image saved to disk instead of volatile memory

# Linux Kernel PM example: CPUidle core

- CPU is running idle
  - Next event (in time) determines how long the CPU can sleep
  - CPU woken up by IRQs, behaves as WFI regardless of the idle state power depth
- Idle state entered through CPUidle driver enter() function
  - Prepare the CPU for power down
  - Execute the power down command

```
1   CPU idle state data
2
3   struct cpuidle_state {
4         char name[CPUIDLE_NAME_LEN];
5         char desc[CPUIDLE_DESC_LEN];
6         unsigned int flags;
7         unsigned int exit_latency; /* in US */
8         int power_usage; /* in mW */
9         unsigned int target_residency; /* in US */
10        bool disabled; /* disabled on all CPUs */
11        /* Idle state enter function */
12        {int (*enter)(struct cpuidle_device *dev,
13                      struct cpuidle_driver *drv,
14                      int index);
15        ...
16  };
```

**ARM**

# ARM 32-bit CPUidle consolidation issues

- Current ARM CPUidle drivers in the kernel cannot be easily cleaned-up/consolidated
- ARM 32-bit SMP CPUidle code was merged separately through different trees, with no consolidation effort
- `cpu_suspend()` kernel interface was defined, but there are still open issues that prevent consolidation
  - Lack of standard power down procedures
  - Lack of standard idle state entry interface
- Multi-Cluster Power Management (MCPM) introduction helped define a common entry point for idle
  - It requires the kernel to run in secure world unless security is overriden in firmware
  - It relies on early devices initialization in the kernel

**ARM**

# ARM 32-bit CPUidle consolidation issues

- Current ARM CPUidle drivers in the kernel cannot be easily cleaned-up/consolidated
- ARM 32-bit SMP CPUidle code was merged separately through different trees, with no consolidation effort
- `cpu_suspend()` kernel interface was defined, but there are still open issues that prevent consolidation
    - Lack of standard power down procedures
    - Lack of standard idle state entry interface
- Multi-Cluster Power Management (MCPM) introduction helped define a common entry point for idle
    - It requires the kernel to run in secure world unless security is overriden in firmware
    - It relies on early devices initialization in the kernel

**ARM**

# ARM 32-bit v7 shutdown procedure (1/3)

```
1   ENTRY(tegra_disable_clean_inv_dcache)
2           stmfd    sp!, {r0, r4-r5, r7, r9-r11, lr}
3           dmb                                       @ ensure ordering
4
5           /* Disable the D-cache */
6           mrc      p15, 0, r2, c1, c0, 0
7           bic      r2, r2, #CR_C
8           mcr      p15, 0, r2, c1, c0, 0
9           isb
10
11          /* Flush the D-cache */
12          cmp      r0, #TEGRA_FLUSH_CACHE_ALL
13          blne     v7_flush_dcache_louis
14          bleq     v7_flush_dcache_all
15
16          /* Turn off coherency */
17          exit_smp r4, r5
18
19          ldmfd    sp!, {r0, r4-r5, r7, r9-r11, pc}
20  ENDPROC(tegra_disable_clean_inv_dcache)
```

https://git.kernel.org/cgit/linux/kernel/git/torvalds/linux.git/tree/
arch/arm/mach-tegra/sleep.S?id=refs/tags/v3.16

ARM

# ARM 32-bit v7 shutdown procedure (2/3)

```
1   ENTRY(omap4_finish_suspend)
2           stmfd   sp!, {r4-r12, lr}
3           cmp     r0, #0x0
4           beq     do_WFI                          @ No lowpower state, jump to WFI
5           [...]
6   skip_secure_l1_clean:
7           bl      v7_flush_dcache_all
8
9           mrc     p15, 0, r0, c1, c0, 0
10          bic     r0, r0, #(1 << 2)               @Disable the C bit
11          mcr     p15, 0, r0, c1, c0, 0
12          isb
13
14          bl      v7_flush_dcache_all
15
16          bl      omap4_get_sar_ram_base
17          mov     r8, r0
18          ldr     r9, [r8, #OMAP_TYPE_OFFSET]
19          cmp     r9, #0x1                        @ Check for HS device
20          bne     scu_gp_set
21          mrc     p15, 0, r0, c0, c0, 5           @ Read MPIDR
22          ands    r0, r0, #0x0f
23          ldreq   r0, [r8, #SCU_OFFSET0]
24          ldrne   r0, [r8, #SCU_OFFSET1]
25          mov     r1, #0x00
26          stmfd   r13!, {r4-r12, r14}
27          ldr     r12, =OMAP4_MON_SCU_PWR_INDEX
28          DO_SMC
29          [...]
30  ENDPROC(omap4_finish_suspend)
```

https://git.kernel.org/cgit/linux/kernel/git/torvalds/linux.git/tree/
arch/arm/mach-omap2/sleep44xx.S?id=refs/tags/v3.16

ARM

# ARM 32-bit v7 shutdown procedure (3/3)

```
1  #define v7_exit_coherency_flush(level) \
2          asm volatile( \
3          "stmfd          sp!, {fp, ip} \n\t" \
4          "mrc            p15, 0, r0, c1, c0, 0       @ get SCTLR \n\t" \
5          "bic            r0, r0, #"__stringify(CR_C)" \n\t" \
6          "mcr            p15, 0, r0, c1, c0, 0       @ set SCTLR \n\t" \
7          "isb            \n\t" \
8          "bl             v7_flush_dcache_"__stringify(level)" \n\t" \
9          "clrex          \n\t" \
10         "mrc            p15, 0, r0, c1, c0, 1       @ get ACTLR \n\t" \
11         "bic            r0, r0, #(1 << 6)       @ disable local coherency \n\t" \
12         "mcr            p15, 0, r0, c1, c0, 1       @ set ACTLR \n\t" \
13         "isb            \n\t" \
14         "dsb            \n\t" \
15         "ldmfd          sp!, {fp, ip}" \
16         : : : "r0","r1","r2","r3","r4","r5","r6","r7", \
17              "r9","r10","lr","memory" )
```

https://git.kernel.org/cgit/linux/kernel/git/torvalds/linux.git/tree/
arch/arm/include/asm/cacheflush.h?id=refs/tags/v3.16

ARM

# ARM 32-bit CPUidle consolidation issues

- Current ARM CPUidle drivers in the kernel cannot be easily cleaned-up/consolidated
- ARM SMP CPUidle code merged separately through different trees, with no consolidation effort
- `cpu_suspend()` kernel interface was defined, but there are still open issues that prevent consolidation
  - Lack of standard power down procedures
  - Lack of standard idle state entry interface
- Multi-Cluster Power Management (MCPM) introduction helped define a common entry point for idle
  - It requires the kernel to run in secure world unless security is overriden in firmware
  - It relies on early devices initialization in the kernel

**ARM**

# ARM 32-bit CPUidle enter examples (1/3)

```
 1  static int armada_370_xp_enter_idle(struct cpuidle_device *dev,
 2                                      struct cpuidle_driver *drv,
 3                                      int index)
 4  {
 5          int ret;
 6          bool deepidle = false;
 7          cpu_pm_enter();
 8
 9          if (drv->states[index].flags & ARMADA_370_XP_FLAG_DEEP_IDLE)
10                  deepidle = true;
11
12          ret = armada_370_xp_cpu_suspend(deepidle);
13          if (ret)
14                  return ret;
15
16          cpu_pm_exit();
17
18          return index;
19  }
```

**ARM**

# ARM 32-bit CPUidle enter examples (2/3)

```
1   static int idle_finisher(unsigned long flags)
2   {
3           exynos_enter_aftr();
4           cpu_do_idle();
5
6           return 1;
7   }
8
9   static int exynos_enter_core0_aftr(struct cpuidle_device *dev,
10                                     struct cpuidle_driver *drv,
11                                     int index)
12  {
13          cpu_pm_enter();
14          cpu_suspend(0, idle_finisher);
15          cpu_pm_exit();
16
17          return index;
18  }
19
20  static int exynos_enter_lowpower(struct cpuidle_device *dev,
21                                   struct cpuidle_driver *drv,
22                                   int index)
23  {
24          int new_index = index;
25
26          if (num_online_cpus() > 1 || dev->cpu != 0)
27                  new_index = drv->safe_state_index;
28
29          if (new_index == 0)
30                  return arm_cpuidle_simple_enter(dev, drv, new_index);
31          else
32                  return exynos_enter_core0_aftr(dev, drv, new_index);
33  }
```
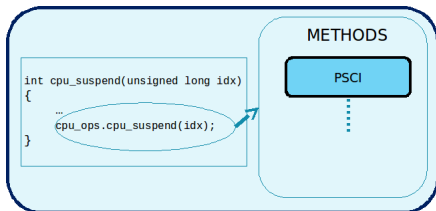
ARM

# ARM 32-bit CPUidle enter examples (3/3)

```
1   static int notrace bl_powerdown_finisher(unsigned long arg)
2   {
3           /* MCPM works with HW CPU identifiers */
4           unsigned int mpidr = read_cpuid_mpidr();
5           unsigned int cluster = MPIDR_AFFINITY_LEVEL(mpidr, 1);
6           unsigned int cpu = MPIDR_AFFINITY_LEVEL(mpidr, 0);
7
8           mcpm_set_entry_vector(cpu, cluster, cpu_resume);
9           [...]
10
11          mcpm_cpu_suspend(0);
12
13          /* return value != 0 means failure */
14          return 1;
15  }
16
17  static int bl_enter_powerdown(struct cpuidle_device *dev,
18                                struct cpuidle_driver *drv, int idx)
19  {
20          cpu_pm_enter();
21
22          cpu_suspend(0, bl_powerdown_finisher);
23
24          mcpm_cpu_powered_up();
25
26          cpu_pm_exit();
27
28          return idx;
29  }
```

https://git.kernel.org/cgit/linux/kernel/git/torvalds/linux.git/tree/
drivers/cpuidle/cpuidle-big_little.c?id=refs/tags/v3.16

ARM

# Linux Kernel ARM64 CPU operations

- Generic kernel interface for CPU operations
  - boot
  - hotplug
  - idle

- CPU operations represent the interface through which the kernel carries out the required actions on CPUs
  - Interface hides the actual method implementation
  - Hooks initialized at boot through DT or ACPI



```
 1  struct cpu_operations {
 2          const char *name;
 3          int (*cpu_init)(struct device_node *,
 4                          unsigned int);
 5          int (*cpu_init_idle)(struct device_node *,
 6                          unsigned int);
 7          int (*cpu_prepare)(unsigned int);
 8          int (*cpu_boot)(unsigned int);
 9          void (*cpu_postboot)(void);
10          int (*cpu_disable)(unsigned int cpu);
11          void (*cpu_die)(unsigned int cpu);
12          int (*cpu_suspend)(unsigned long);
13  };
```

ARM

# ARM64 CPUidle enter

```
1  static int arm_enter_idle_state(struct cpuidle_device *dev,
2                                  struct cpuidle_driver *drv, int idx)
3  {
4          int ret;
5
6          if (!idx) {
7                  cpu_do_idle();
8                  return idx;
9          }
10
11         ret = cpu_pm_enter();
12         if (!ret) {
13                 /*
14                  * Pass idle state index to cpu_suspend which in turn will
15                  * call the CPU ops suspend protocol with idle index as a
16                  * parameter.
17                  */
18                 ret = cpu_suspend(idx);
19
20                 cpu_pm_exit();
21         }
22
23         return ret ? -1 : idx;
24  }
```

http://lists.infradead.org/pipermail/linux-arm-kernel/2014-July/
274249.html

# ARM64 CPU Operations: PSCI suspend

```
1  static int __maybe_unused cpu_psci_cpu_suspend(unsigned long index)
2  {
3          struct psci_power_state *state = __get_cpu_var(psci_power_state);
4          /*
5           * idle state index 0 corresponds to wfi, should never be called
6           * from the cpu_suspend operations
7           */
8          if (WARN_ON_ONCE(!index))
9                  return -EINVAL;
10
11         return psci_ops.cpu_suspend(state[index - 1],
12                                 virt_to_phys(cpu_resume));
13 }
14
```

http:
//lists.infradead.org/pipermail/linux-arm-kernel/2014-July/274245.html

- Suspend operations become a stub that calls into firmware to carry out power down
  - Equivalent of `mwait` in x86 world

ARM

# Idle States Device Tree Bindings

```
1  idle-states {
2          entry-method = "arm,psci";
3          CPU_SLEEP_0: cpu-sleep-0 {
4                  compatible = "arm,idle-state";
5                  arm,psci-suspend-param = <0x0010000>;
6                  entry-latency-us = <40>;
7                  exit-latency-us = <100>;
8                  min-residency-us = <150>;
9          };
10         CLUSTER_SLEEP_0: cluster-sleep-0 {
11                 compatible = "arm,idle-state";
12                 arm,psci-suspend-param = <0x1010000>;
13                 entry-latency-us = <500>;
14                 exit-latency-us = <1000>;
15                 min-residency-us = <2500>;
16         };
17 };
```

http://lists.infradead.org/pipermail/linux-arm-kernel/2014-July/274251.html

- Configuration data provided by firmware through well established bindings

ARM

# ARM64 Kernel PM: What's to come

- Augment DT idle states bindings with power domains information
  - Link CPU components (inclusive of caches) and devices to their respective power domains
  - Add power domains information to DT idle states
- Implement ACPI ARM PM drivers for ACPI 5.1
- Spec out ARM system sleep states
  - PSCI system sleep states management
  - ACPI and DT system sleep states bindings

**ARM**

# Conclusion

- ARM 32-bit Linux kernel huge legacy code
    - HW power control design quirks and related SW workarounds
    - Lack of power management standards in both HW and SW
    - Lots of out-of-tree power management ARM vendors code
- ARM 64-bit power management
    - No legacy
    - ARM 32-bit SW design experience
    - Upstream momentum for standard interfaces and systems
- Standardization required in both HW and SW to prevent same ARM 32-bit mistakes
    - HW power management recommendation to be shared with ARM vendors
    - Foster PSCI standardization deployment
    - Actively contribute to ACPI specifications and device tree bindings

**ARM**

# THANKS !!!