# Knocking at your back door (O.H.D.W.M.I.A.C.A.Y.S.)

**ARM**

Marc Zyngier `<marc.zyngier@arm.com>`

ELCE16

October 13, 2016

# Content

- Basics of an interrupt
- Interrupt controllers
- Linux's data structures
- Chained interrupt controllers
- Hierarchical interrupt controllers
- Generic MSIs
- ...
- Profit!

**ARM**

# Please interrupt me

**ARM**

# Talk you should not have missed

- IRQs: the Hard, the Soft, the Threaded and the Preemptible
- Alison Chaiken, Peloton Technology
- Took place on Tuesday[1]
- Covers the dynamic aspects of interrupt handling

---

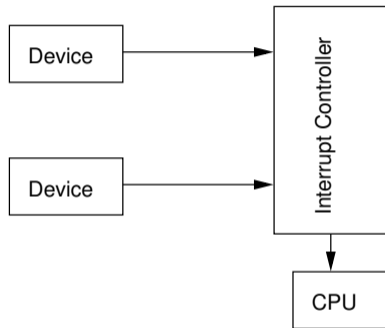[1]Use your TARDIS or wait for it to appear on some website

**ARM**

# What is an interrupt?

- A hardware signal
- Emited from a peripheral to a CPU
- Indicating that a device-specific condition has been satisfied
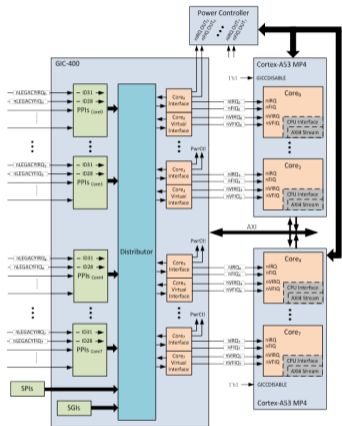
**ARM**

# Multiplexing interrupts

- Having a single interrupt for the CPU is usually not enough
- Most systems have tens, hundreds of them
- An interrupt controller allows them to be multiplexed
- Very often architecture or platform specific
- Offers specific facilities
  - Masking/unmasking individual interrupts
  - Setting priorities
  - SMP affinity
  - Exotic things like wake-up interrupts

**ARM**

# Multiplexing interrupts

- Having a single interrupt for the CPU is usually not enough

- Most systems have tens, hundreds of them

- An interrupt controller allows them to be multiplexed

- Very often architecture or platform specific

- Offers specific facilities

  - Masking/unmasking individual interrupts
  - Setting priorities
  - SMP affinity
  - Exotic things like wake-up interrupts



GIC-400, simplified view

**ARM**

# Interrupt triggers

- Level triggered (high or low)
  - Indicates a persistent condition
  - An action has to be performed on the device to clear the interrupt

- Edge triggered (rising or falling)
  - Indicates an event
  - May have happened once or more...

- Some systems do not expose the trigger type to software
  - Either the interrupt is abstracted (virtualization)
  - Or this is more an exception than an interrupt...

**ARM**

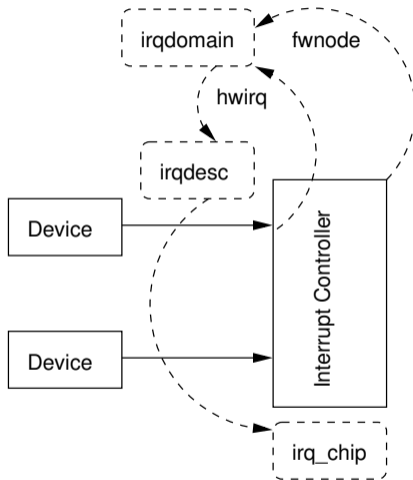# "And now for something completely different..."

Monty Python's Flying Circus

**ARM**

# How does Linux deal with interrupts

- `struct irq_chip`
  - A set of methods describing how to drive the interrupt controller
  - Directly called by core IRQ code
- `struct irqdomain`
  - A pointer to the firmware node for a given interrupt controller (`fwnode`)
  - A method to convert a firmware description of an IRQ into an ID local to this interrupt controller (`hwirq`)
  - A way to retrieve the Linux view of an IRQ from the `hwirq`
- `struct irq_desc`
  - Linux's view of an interrupt
  - Contains all the core stuff
  - 1:1 mapping to the Linux interrupt number
- `struct irq_data`
  - Contains the data that is relevant to the `irq_chip` managing this interrupt
    - Both the Linux IRQ number and the `hwirq`
    - A pointer to the `irq_chip`
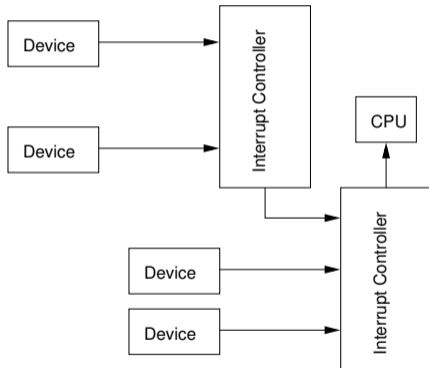    - Embedded in `irq_desc` (for now)

**ARM**

# In a nutshell

- CPU gets an interrupt
- Find out the `hwirq` from the interrupt controller
  - Usually involves reading some HW register
- Look-up the `irq_desc` into the `irqdomain` using the `hwirq`
  - Actually returns an IRQ number, which is equivalent to the `irq_desc`
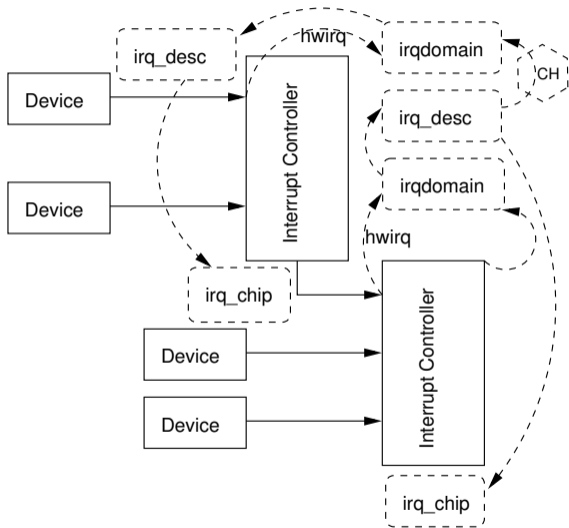- The core kernel then handles the interrupt

**ARM**

# Multiplexing more interrupts

- Not enough interrupts lines?
  - Dedicate a single line for a secondary interrupt controller
  - And add more stuff to it!

- Requires two level handling
  - First handle the interrupt on the primary interrupt controller
  - Then at the secondary one to find out which device has caused the interrupt
  - See `irq_set_chained_handler_and_data`, `chained_irq_enter`, `chained_irq_exit`
  - Never treat this as a normal interrupt handler

- Used in each and every x86 system
  - The infamous i8259 cascade

- You can also share a single interrupt between devices
  - And that really stinks. Please avoid doing it if possible.

ARM

# Chained irqchips, the irqdomain view

- Each interrupt controller has its own `irqdomain`
- The kernel deals with two interrupts
  - and two interrupt handlers
  - the first one being a `chained` handler
  - convention is to stash a pointer to the secondary domain inside the top-level `irq_desc`
- We walk the interrupt chain in reverse order
- Once we reach the last level `irq_desc`, we can process the actual interrupt handler
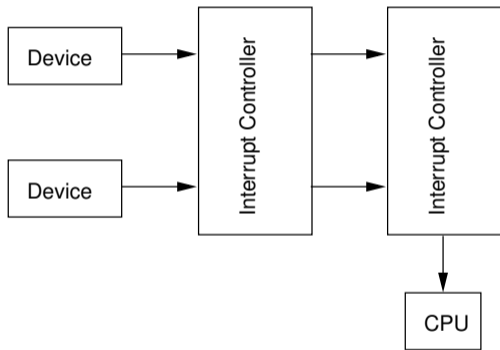
**ARM**

# The DT view

- A secondary irqchip points to the one implementing the first level
  - Use `interrupts` to describe the signal path between irqchips
  - The secondary chip owns the cascade interrupt
  - It doesn't appear in `/proc/interrupts`
- Use `interrupt-parent` to point the device at the right interrupt controller

```
1  interrupt-parent = <&gic>;
2
3  gic: interrupt-controller@01c81000 {
4          compatible = "arm,cortex-a7-gic", "arm,cortex-a15-gic";
5          interrupt-controller;
6          #interrupt-cells = <3>;
7          interrupts = <GIC_PPI 9 (GIC_CPU_MASK_SIMPLE(4) |
8                                   IRQ_TYPE_LEVEL_HIGH)>;
9  };
10
11 nmi_intc: interrupt-controller@01c00030 {
12          compatible = "allwinner,sun7i-a20-sc-nmi";
13          interrupt-controller;
14          #interrupt-cells = <2>;
15          interrupts = <GIC_SPI 0 IRQ_TYPE_LEVEL_HIGH>;
16 };
17
18 axp209: pmic@34 {
19          interrupt-parent = <&nmi_intc>;
20          interrupts = <0 IRQ_TYPE_LEVEL_LOW>;
21 };
```
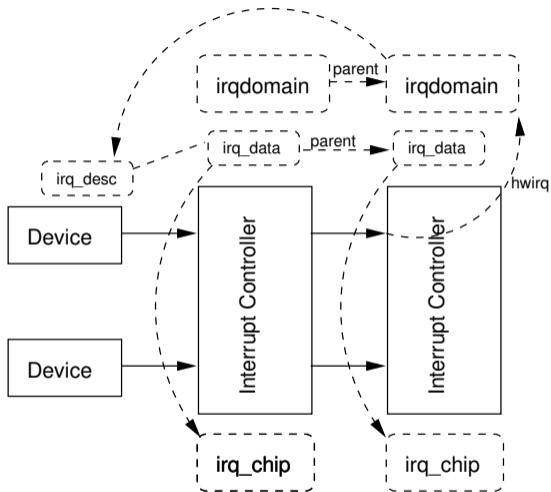
**ARM**

# When multiplexing doesn't fit

- There is more than just cascading irqchips
- Some setups have a 1:1 mapping between input and output
  - Interrupt routers
  - Wake-up controllers
  - Programmable line inverters
- Most of them are not interrupt controllers
  - Still, they do impact the interrupt delivery
  - We choose to represent them as `irq_chip`
- This is a hierarchical/stacked configuration
- The chained irqchip paradigm doesn't match it

**ARM**

# Hierarchical (stacked) IRQ domains

- We want the same `irq_desc` to be valid across all irqchips
  - This ensures that the Linux IRQ number is unique for a given signal path

- For a given `irq_desc`, each irqchip should be responsible for the `hwirq`
  - This fits the `irq_data` properties

- Most of the data structures now have a `parent` field representing the hierarchy

- The handling is done by walking the signal path in delivery order
  - A given irqchip can perform some local action before forwarding the request to its parent
  - Or even terminate the handling early



© ARM 2016

**ARM**

# Hierarchical domains, the DT view

- Each intermediate irqchip points to its parent
  - Do not use `interrupts` to describe the signal path between irqchips
  - Use a device-specific property to decribe an interrupt range/space if necessary
- The root irqchip points to itself
  - A DT oddity...
- Devices can point to any element of the stack
  - The device interrupt specifiers must match the first irqchip in the signal path

```
1  interrupt-parent = <&sysirq>;
2
3  sysirq: intpol-controller@10200620 {
4          interrupt-controller;
5          #interrupt-cells = <3>;
6          interrupt-parent = <&gic>;
7  };
8
9  gic: interrupt-controller@10231000 {
10         #interrupt-cells = <3>;
11         interrupt-parent = <&gic>;
12         interrupt-controller;
13 };
14
15 uart0: serial@11002000 {
16         interrupts = <GIC_SPI 91 IRQ_TYPE_LEVEL_LOW>;
17 };
```

**ARM**

# "Message in a bottle"

The Police, Reggatta de Blanc

**ARM**

# More than wired interrupts: MSIs

Message Signaled Interrupts are an essential part of the interrupt infrastructure

- A simple 32bit write (the message) from the device to a doorbell
    - The doorbell is usually the interrupt controller itself
    - The generated interrupt depends on the data being written
    - By definition edge triggered
- Avoid the spider web syndrome
    - Routing interrupts to the periphery of a SoC is a constraint
    - MSIs allows the use of the same busses as the data
    - Having multiple interrupts per device costs nothing
- Acts as a memory barrier w.r.t DMA
    - Avoid the "got an interrupt but data is not there yet" problem
- Bus agnostic
    - Historically tied to PCI(e)
    - Now implemented on all kinds of busses...

© ARM 2016                                                                                          **ARM**

# The goals of supporting MSIs in a generic way

- We'd like to support MSIs on any bus
- We want to cater for the weird and wonderful stuff
  - Intel's DMAR
  - ARM's GICv3 ITS
  - Freescale's MC bus
  - Platform devices
  - Hisilicon's MBIGEN
- Must nicely cohabit with the current PCI/MSI implementation
- Hierarchical domains are a good solution for this[2]
  - Entirely implemented as part of the core IRQ code (`kernel/irq/msi.c`)
  - Per-bus front-ends
    - `drivers/pci/msi.c`
    - `drivers/base/platform-msi.c`
    - `drivers/staging/fsl-mc/bus/mc-msi.c`

---

[2]Please trust me on that one...

**ARM**

# Generic MSI

- **irq_chip grows two new methods**
  - **irq_compose_msi_msg**: populate a msi_msg
    - Address of the doorbell + data to be written
    - Implemented by the MSI controller, bus agnostic
  - **irq_write_msi_msg**
    - Write the content of the msi_msg to a given device
    - Implemented by the bus front-end, bus specific

- **msi_domain_info to describe a MSI domain**
  - A struct irq_chip
    - Must at least contain a irq_write_msi_msg method
  - A struct msi_domain_ops
    - A set of functions used to build an irqdomain
  - A set of flags (some bus specific), and allowing most of the above to get sensible defaults

- **Bus specific irqdomain creation functions**

```
1  /*
2   * PCI/MSI setup
3   */
4  static struct irq_chip my_msi_irq_chip = {
5          .name            = "MSI",
6          .irq_eoi         = irq_chip_eoi_parent,
7          .irq_write_msi_msg = pci_msi_domain_write_msg,
8  };
9
10 static struct msi_domain_info my_msi_dom_info = {
11         .flags       = (MSI_FLAG_USE_DEF_DOM_OPS |
12                         MSI_FLAG_USE_DEF_CHIP_OPS |
13                         MSI_FLAG_PCI_MSIX),
14         .chip        = &my_msi_irq_chip,
15 };
16
17 [...]
18 /*
19  * Build the PCI/MSI domain on top of the IRQ domain
20  * representing the MSI hardware
21  */
22 pci_domain = pci_msi_create_irq_domain(fwnode,
23                                        &my_msi_dom_info,
24                                        irq_domain);
```
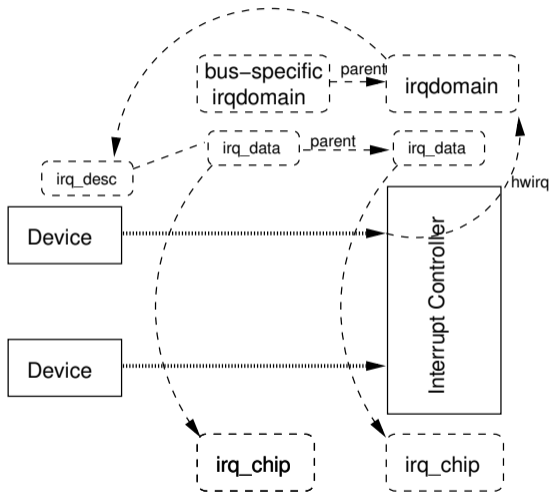
**ARM**

# Generic MSI

- `irq_chip` **grows two new methods**
  - `irq_compose_msi_msg`: **populate a** `msi_msg`
    - Address of the doorbell + data to be written
    - Implemented by the MSI controller, bus agnostic
  - `irq_write_msi_msg`
    - Write the content of the `msi_msg` to a given device
    - Implemented by the bus front-end, bus specific

- `msi_domain_info` **to describe a MSI domain**
  - A `struct irq_chip`
    - Must at least contain a `irq_write_msi_msg` method
  - A `struct msi_domain_ops`
    - A set of functions used to build an `irqdomain`
  - A set of flags (some bus specific), and allowing most of the above to get sensible defaults

- **Bus specific** `irqdomain` **creation functions**

```
1   /*
2    * platform-msi setup
3    */
4   static struct irq_chip my_pmsi_irq_chip = {
5           .name   = "pMSI",
6   };
7
8   static struct msi_domain_ops my_pmsi_ops = {
9   };
10
11  static struct msi_domain_info my_pmsi_dom_info = {
12          .flags  = (MSI_FLAG_USE_DEF_DOM_OPS |
13                     MSI_FLAG_USE_DEF_CHIP_OPS),
14          .ops    = &my_pmsi_ops,
15          .chip   = &my_pmsi_irq_chip,
16  };
17
18  [...]
19  /*
20   * Build the platform-msi domain on top of the IRQ domain
21   * representing the MSI hardware
22   */
23  plat_domain = platform_msi_create_irq_domain(fwnode,
24                                  &my_pmsi_dom_info,
25                                  irq_domain);
```

**ARM**

# Generic MSI in pictures

- **At configuration time**
  - The MSI controller irqchip composes the message
  - The bus-specific irqchip programs the device

- **Everything is just like the stacked irqchip scenario**
  - The only notable difference is that we have a bus-specific irqdomain that doesn't correspond to any HW
  - Its main function is to cater for different programing interfaces at the device level



© ARM 2016

**ARM**

# A platform MSI special

- There is no such thing as a "standard" platform device
- No way to implement a `irq_write_msi_msg` in a standard way
- Worked around by providing it at allocation time
  - The function is per-device
  - Allows for any crazy stuff

```
1  static void arm_smmu_write_msi_msg(struct msi_desc *desc,
2                                     struct msi_msg *msg)
3  {
4    doorbell = (((u64)msg->address_hi) << 32) | msg->address_lo;
5
6    writeq_relaxed(doorbell, smmu->base + cfg[0]);
7    writel_relaxed(msg->data, smmu->base + cfg[1]);
8  }
9
10 static void arm_smmu_setup_msis(struct arm_smmu_device *smmu)
11 {
12   [...]
13   ret = platform_msi_domain_alloc_irqs(dev, nvec,
14                                        arm_smmu_write_msi_msg);
15   [...]
16   for_each_msi_entry(desc, dev) {
17           switch (desc->platform.msi_index) {
18           /* request desc->irq */
19           }
20   }
21 }
```

**ARM**

# "I'm going slightly mad"

Queen, Innuendo

**ARM**

# The interrupt strikes back

- Just as we thought we had fixed the world by giving MSIs to everyone...
- People now build wired interrupt controllers...
- ... that use MSI as their transport
  - Allows wired devices to be placed far away from the irqchip
  - Conveniently, one MSI per wire
- Stacked domains to the rescue!
  - The irqchip is a MSI-capable device
  - We can give it its own irqdomain

**ARM**

# Wire-MSI bridges, the programatic view

- At probe time, create a device-specific domain
- Automatically attached to the device's `msi-parent`'s own domain
- When allocating its MSIs, place them in that domain
- Dish out wired interrupts as a normal irqchip

```
1  static struct irq_domain_ops mbigen_domain_ops = {
2    [...]
3  };
4
5  static int mbigen_irq_domain_alloc(struct irq_domain *domain,
6                                     unsigned int virq,
7                                     unsigned int nr_irqs,
8                                     void *args)
9  {
10   struct irq_fwspec *fwspec = args;
11
12   mbigen_domain_translate(domain, fwspec, &hwirq, &type);
13   platform_msi_domain_alloc(domain, virq, nr_irqs);
14   mgn_chip = platform_msi_get_host_data(domain);
15
16   for (i = 0; i < nr_irqs; i++)
17     irq_domain_set_hwirq_and_chip(domain, virq + i, hwirq + i,
18                                   &mbigen_irq_chip, mgn_chip->base);
19 }
20
21 static struct irq_domain_ops mbigen_domain_ops = {
22         .alloc = mbigen_irq_domain_alloc,
23 };
24
25 static int mbigen_device_probe(struct platform_device *pdev)
26 {
27 [...]
28   domain = platform_msi_create_device_domain(&child->dev,
29                                              num_pins,
30                                              mbigen_write_msg,
31                                              &mbigen_domain_ops,
32                                              mgn_chip);
33 }
```
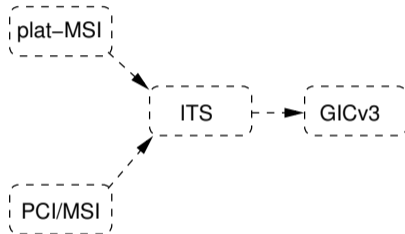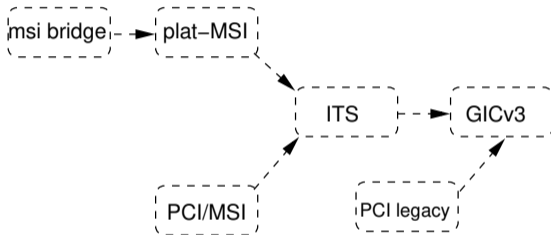
© ARM 2016

**ARM**

# IRQ domains, as seen on `arm64`...

**ARM**

GICv3

**ARM**

# IRQ domains, as seen on `arm64`...

**ARM**

# IRQ domains, as seen on `arm64`...

**ARM**

# IRQ domains, as seen on `arm64`...

**ARM**

**ARM**

# IRQ domains, as seen on `arm64`...

**ARM**

# Thank you!

**ARM**