# Interrupt filtering

Paolo Pisati

Universita' Statale di Milano
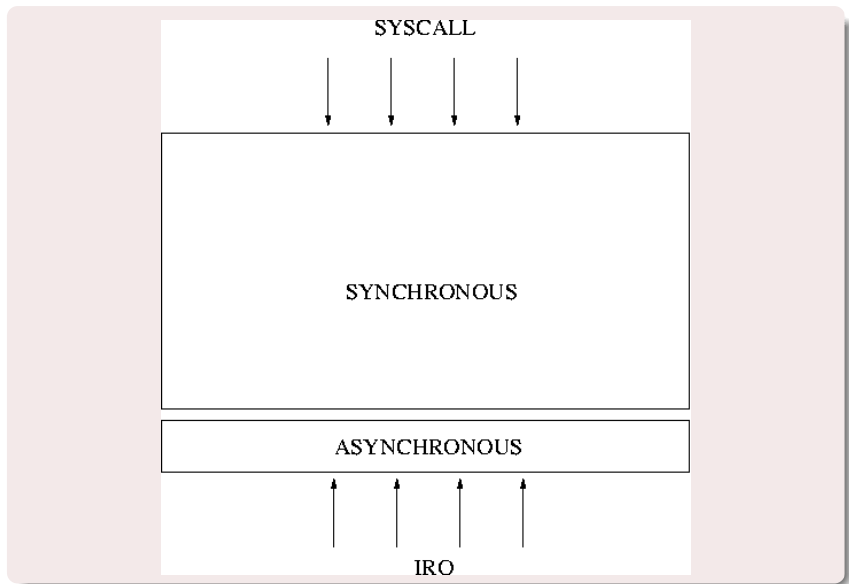
December 1, 2006

### What's this talk about?

- Interrupt handling mechanism

- An interrupt is an asynchronous signal from hardware indicating the need for attention

- Key aspect: Synchronization policy

### What's this talk about?

- Interrupt handling mechanism
- An interrupt is an asynchronous signal from hardware indicating the need for attention
- Key aspect: Synchronization policy

### What's this talk about?

- Interrupt handling mechanism
- An interrupt is an asynchronous signal from hardware indicating the need for attention
- Key aspect: Synchronization policy

## **Outline**

**1** **FreeBSD 4x**

**2** **FreeBSD SMPng**

**3** **Interrupt filtering**

**4** **Performance**
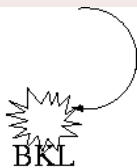
**5** **Implementation details**

## Outline

**1  FreeBSD 4x**

**2**  FreeBSD SMPng

**3**  Interrupt filtering
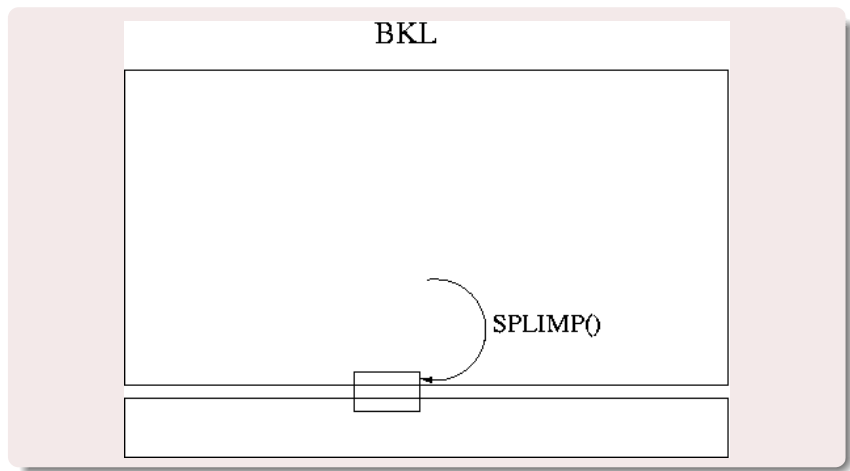
**4**  Performance
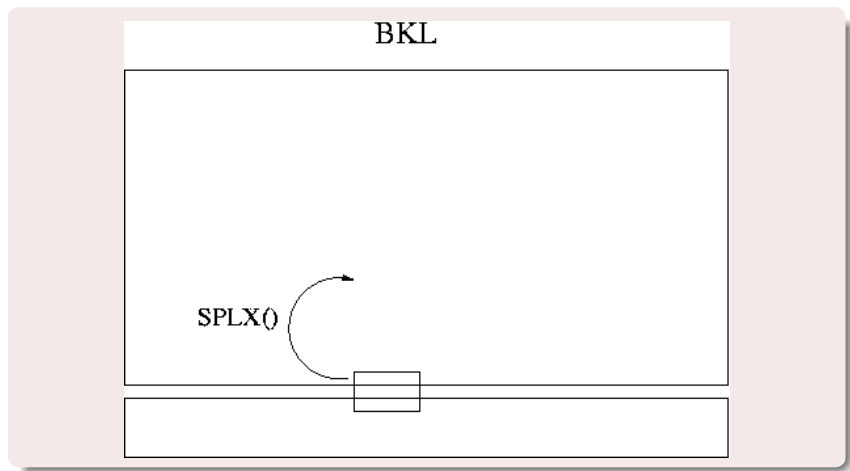
**5**  Implementation details

- Monothread kernel
- BKL ruled access to kernel space
- SPL calls used to synchronize top and bottom halves

- Monothread kernel
- BKL ruled access to kernel space
- SPL calls used to synchronize top and bottom halves

BKL

- Monothread kernel
- BKL ruled access to kernel space
- SPL calls used to synchronize top and bottom halves

BKL

SPLIMP()

# Outline

- Multithreaded kernel
- Interrupt handler got a private context: ithread
- Time critical handlers that don't block, can run in the context of the interrupted process (FAST)

- Multithreaded kernel
- Interrupt handler got a private context: ithread
- Time critical handlers that don't block, can run in the context of the interrupted process (FAST)

MTX_UNLOCK()

MTX

MTX_LOCK()

- Multithreaded kernel
- Interrupt handler got a private context: ithread
- Time critical handlers that don't block, can run in the context of the interrupted process (FAST)

## The problems

- Interrupt latency
- ...even worse with shared irq
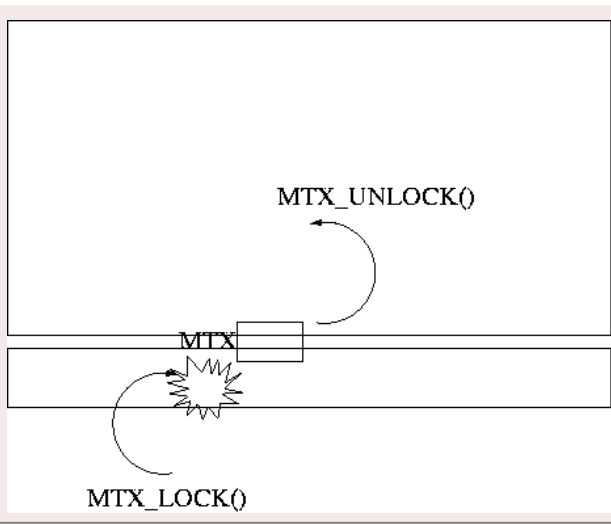- Interrupt masking at controller level

### The problems

- Interrupt latency
- ...even worse with shared irq
- Interrupt masking at controller level

### The problems

- Interrupt latency
- ...even worse with shared irq
- Interrupt masking at controller level

### The problems

- Interrupt latency
- ...even worse with shared irq
- Interrupt masking at controller level

## Outline

- Divide the interrupt handler in 2 logical pieces:

  **FILTER** runs in interrupt context, checks the received
  interrupt, serves it/delegates more work
  **ITHREAD** runs in ithread context, can block

  - ack the interrupt after filter execution

- Divide the interrupt handler in 2 logical pieces:
  **FILTER** runs in interrupt context, checks the received interrupt, serves it/delegates more work
  **ITHREAD** runs in ithread context, can block

- ack the interrupt after filter execution

- Divide the interrupt handler in 2 logical pieces:
  **FILTER** runs in interrupt context, checks the received
  interrupt, serves it/delegates more work
  **ITHREAD** runs in ithread context, can block

- ack the interrupt after filter execution

- Divide the interrupt handler in 2 logical pieces:
  **FILTER** runs in interrupt context, checks the received interrupt, serves it/delegates more work
  **ITHREAD** runs in ithread context, can block
- ack the interrupt after filter execution

ITHREAD                                                              FILTER+ITHREAD

ISR()                                                                ISR()

    MASK INT()                                                           CALL FILTER()
    ASK ITHREAD SCHED                                                    ACK INT()
                                                                         ASK ITHREAD SCHED

IRET
 ...                        SCHEDULER WAKEUP                         IRET
 ...                                                                  ...
 ...                                                                  ...
 ...                                                                  ...
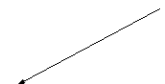ITHREAD0()                                                            ...
 ...                                                                 ITHREAD()
 ...                                                                  ...
RET                                                                  RET
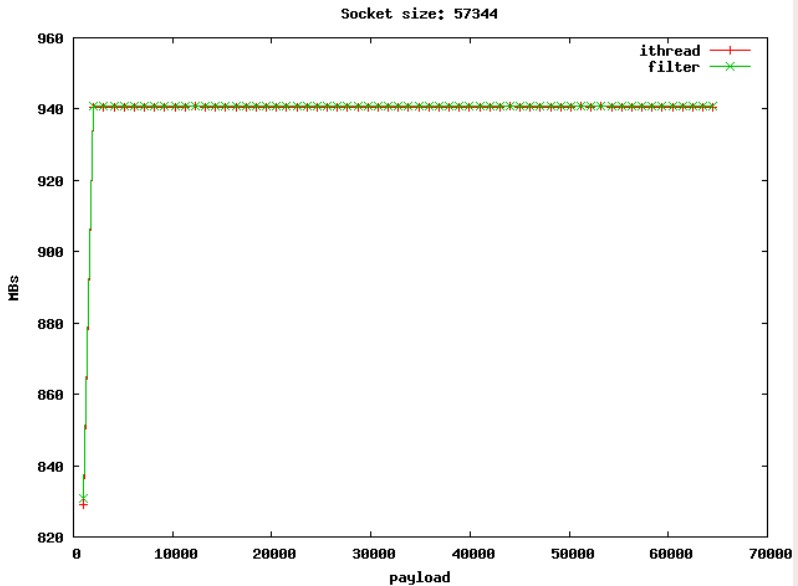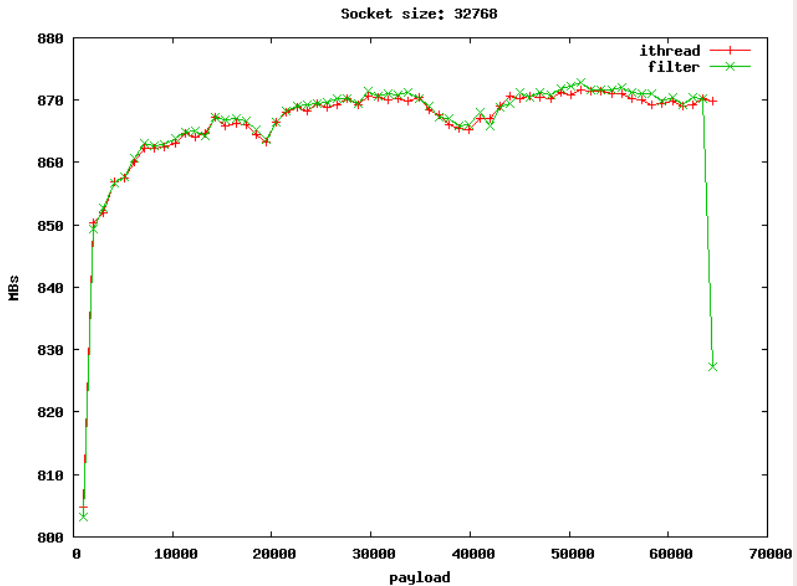 ...
 ...
DEMASK AND ACK INT()

Let's see some code: bfe, em, re, aac, xl

## Outline

1. FreeBSD 4x

2. FreeBSD SMPng

3. Interrupt filtering

4. **Performance**

5. Implementation details

Socket size: 32768

Filter vs Ithread

4.x vs Filter vs Ithread

# Outline

## Implementation details

- the total patch against HEAD is about 300kb
- newbus API change: bus_setup_intr()
- 3 different ways to handle interrupts:

## Implementation details

- the total patch against HEAD is about 300kb
- newbus API change: bus_setup_intr()
- 3 different ways to handle interrupts:

- ```
  int
  bus_setup_intr(
      device_t dev,
      struct resource *r,
      int flags,
      driver_filter_t filter,
      driver_intr_t handler,
      void *arg,
      void **cookiep
  );
  ```

- ```
  int driver_filter_t(void*);
  ```

- int
  bus_setup_intr(
      device_t dev,
      struct resource *r,
      int flags,
      driver_filter_t filter,
      driver_intr_t handler,
      void *arg,
      void **cookiep
  );

- int driver_filter_t(void*);

## Implementation details

- the total patch against HEAD is about 300kb
- newbus API change: bus_setup_intr()
- 3 different ways to handle interrupts:

  **FILTER (aka FAST)** bus_setup_intr(dev, r, flags, filter, NULL, arg, cookiep);

  **ITHREAD** bus_setup_intr(dev, r, flags, NULL, ithread, arg, cookiep);

  **FILTER+ITHREAD** bus_setup_intr(dev, r, flags, filter, ithread, arg, cookiep);

### Implementation details

- the total patch against HEAD is about 300kb
- newbus API change: bus_setup_intr()
- 3 different ways to handle interrupts:

  **FILTER (aka FAST)** bus_setup_intr(dev, r, flags, filter, NULL, arg, cookiep);

  **ITHREAD** bus_setup_intr(dev, r, flags, NULL, ithread, arg, cookiep);

  **FILTER+ITHREAD** bus_setup_intr(dev, r, flags, filter, ithread, arg, cookiep);

### Implementation details

- the total patch against HEAD is about 300kb
- newbus API change: bus_setup_intr()
- 3 different ways to handle interrupts:
  **FILTER (aka FAST)** bus_setup_intr(dev, r, flags, filter, NULL, arg, cookiep);
  **ITHREAD** bus_setup_intr(dev, r, flags, NULL, ithread, arg, cookiep);
  **FILTER+ITHREAD** bus_setup_intr(dev, r, flags, filter, ithread, arg, cookiep);

### Implementation details

- the total patch against HEAD is about 300kb
- newbus API change: bus_setup_intr()
- 3 different ways to handle interrupts:
  **FILTER (aka FAST)** bus_setup_intr(dev, r, flags, filter, NULL, arg, cookiep);
  **ITHREAD** bus_setup_intr(dev, r, flags, NULL, ithread, arg, cookiep);
  **FILTER+ITHREAD** bus_setup_intr(dev, r, flags, filter, ithread, arg, cookiep);

## Implementation details(2)

- filter have 3 returns code:

    **FILTER STRAY** event not recognized
    **FILTER HANDLED** interrupt is acked/turned off
    **FILTER SCHEDULE THREAD** schedule the ithread

- no other value can be returned with FILTER STRAY

- if a filter wants to schedule an ithread, it returns
  FILTER HANDLED | FILTER SCHEDULE THREAD

- if all filters returned FILTER STRAY...

- or no handlers were registered on that line...

- a new interrupt mitigation mechanism will kick in

## Implementation details(2)

- filter have 3 returns code:

  **FILTER_STRAY** event not recognized

  FILTER_HANDLED interrupt is acked/turned off

  FILTER_SCHEDULE_THREAD schedule the ithread

- no other value can be returned with FILTER_STRAY

- if a filter wants to schedule an ithread, it returns
  FILTER_HANDLED | FILTER_SCHEDULE_THREAD

- if all filters returned FILTER_STRAY...

- or no handlers were registered on that line...

- a new interrupt mitigation mechanism will kick in

## Implementation details(2)

- filter have 3 returns code:
  **FILTER_STRAY** event not recognized
  **FILTER_HANDLED** interrupt is acked/turned off
  **FILTER_SCHEDULE_THREAD** schedule the ithread

- no other value can be returned with FILTER_STRAY

- if a filter wants to schedule an ithread, it returns
  FILTER_HANDLED | FILTER_SCHEDULE_THREAD

- if all filters returned FILTER_STRAY...

- or no handlers were registered on that line...

- a new interrupt mitigation mechanism will kick in

## Implementation details(2)

- filter have 3 returns code:

  **FILTER_STRAY** event not recognized
  **FILTER_HANDLED** interrupt is acked/turned off
  **FILTER_SCHEDULE_THREAD** schedule the ithread

- no other value can be returned with FILTER_STRAY

- if a filter wants to schedule an ithread, it returns
  FILTER_HANDLED | FILTER_SCHEDULE_THREAD

- if all filters returned FILTER_STRAY...

- or no handlers were registered on that line...

- a new interrupt mitigation mechanism will kick in

### Implementation details(2)

- filter have 3 returns code:
  **FILTER STRAY** event not recognized
  **FILTER HANDLED** interrupt is acked/turned off
  **FILTER SCHEDULE THREAD** schedule the ithread

- no other value can be returned with FILTER STRAY

- if a filter wants to schedule an ithread, it returns
  FILTER HANDLED | FILTER SCHEDULE THREAD

- if all filters returned FILTER STRAY...

- or no handlers were registered on that line...

- a new interrupt mitigation mechanism will kick in

### Implementation details(2)

- filter have 3 returns code:

  **FILTER_STRAY** event not recognized
  **FILTER_HANDLED** interrupt is acked/turned off
  **FILTER_SCHEDULE_THREAD** schedule the ithread

- no other value can be returned with FILTER_STRAY

- if a filter wants to schedule an ithread, it returns
  FILTER_HANDLED | FILTER_SCHEDULE_THREAD

- if all filters returned FILTER_STRAY...

- or no handlers were registered on that line...

- a new interrupt mitigation mechanism will kick in

### Implementation details(2)

- filter have 3 returns code:

  **FILTER_STRAY** event not recognized
  **FILTER_HANDLED** interrupt is acked/turned off
  **FILTER_SCHEDULE_THREAD** schedule the ithread

- no other value can be returned with FILTER_STRAY

- if a filter wants to schedule an ithread, it returns
  FILTER_HANDLED | FILTER_SCHEDULE_THREAD

- if all filters returned FILTER_STRAY...

- or no handlers were registered on that line...

- a new interrupt mitigation mechanism will kick in

### Implementation details(2)

- filter have 3 returns code:

  **FILTER_STRAY** event not recognized

  **FILTER_HANDLED** interrupt is acked/turned off

  **FILTER_SCHEDULE_THREAD** schedule the ithread

- no other value can be returned with FILTER_STRAY

- if a filter wants to schedule an ithread, it returns FILTER_HANDLED | FILTER_SCHEDULE_THREAD

- if all filters returned FILTER_STRAY...

- or no handlers were registered on that line...

- a new interrupt mitigation mechanism will kick in

### Implementation details(2)

- filter have 3 returns code:
  **FILTER_STRAY** event not recognized
  **FILTER_HANDLED** interrupt is acked/turned off
  **FILTER_SCHEDULE_THREAD** schedule the ithread

- no other value can be returned with FILTER_STRAY
- if a filter wants to schedule an ithread, it returns
  FILTER_HANDLED | FILTER_SCHEDULE_THREAD
- if all filters returned FILTER_STRAY...
- or no handlers were registered on that line...
- a new interrupt mitigation mechanism will kick in

## Implementation details(3)

- interrupt.h::struct intr_event was modified
- kern_intr.c::intr_event_create() was modified too
- int (*ie_pending)(void *);
- PowerPC MD code was modified to accept more than one FAST handler per line(previouslys INTR_FAST implied INTR_EXCL)

### Implementation details(3)

- interrupt.h::struct intr_event was modified
- kern_intr.c::intr_event_create() was modified too
- int (*ie_pending)(void *);
- PowerPC MD code was modified to accept more than one FAST handler per line(previouslys INTR_FAST implied INTR_EXCL)

## Implementation details(3)

- interrupt.h::struct intr_event was modified
- kern_intr.c::intr_event_create() was modified too
- int (*ie_pending)(void *);
- PowerPC MD code was modified to accept more than one FAST handler per line(previouslys INTR_FAST implied INTR_EXCL)

### Implementation details(3)

- interrupt.h::struct intr_event was modified
- kern_intr.c::intr_event_create() was modified too
- int (*ie_pending)(void *);
- PowerPC MD code was modified to accept more than one FAST handler per line(previouslys INTR_FAST implied INTR_EXCL)

### TODO

- Check the MD code that turn off/mask interrupts (arm, ia64, powerpc and sparc64)

- Check drivers that override newbus generic_bus_setup_intr() (dev/puc/puc.c)

- Much of the MD code in intr_execute_handlers() could be turned into MI

- Testing

### TODO

- Check the MD code that turn off/mask interrupts (arm, ia64, powerpc and sparc64)
- Check drivers that override newbus generic_bus_setup_intr() (dev/puc/puc.c)
- Much of the MD code in intr_execute_handlers() could be turned into MI
- Testing

### TODO

- Check the MD code that turn off/mask interrupts (arm, ia64, powerpc and sparc64)
- Check drivers that override newbus generic_bus_setup_intr() (dev/puc/puc.c)
- Much of the MD code in intr_execute_handlers() could be turned into MI
- Testing

### TODO

- Check the MD code that turn off/mask interrupts (arm, ia64, powerpc and sparc64)
- Check drivers that override newbus generic_bus_setup_intr() (dev/puc/puc.c)
- Much of the MD code in intr_execute_handlers() could be turned into MI
- Testing

### TODO(2)

- Do more measurements
- Develop some tools to measure interrupt latency
- Take a look at lock-free/wait-free data struct

## TODO(2)

- Do more measurements
- Develop some tools to measure interrupt latency
- Take a look at lock-free/wait-free data struct

### TODO(2)

- Do more measurements
- Develop some tools to measure interrupt latency
- Take a look at lock-free/wait-free data struct

## Conclusion

- No performance regression against previous model

- Lower interrupt latency

- Less interrupt problems ("storms")

- I wish to make it part of 7.x

### Conclusion

- No performance regression against previous model
- Lower interrupt latency
- Less interrupt problems ("storms")
- I wish to make it part of 7.x

### Conclusion

- No performance regression against previous model
- Lower interrupt latency
- Less interrupt problems ("storms")
- I wish to make it part of 7.x

### Conclusion

- No performance regression against previous model
- Lower interrupt latency
- Less interrupt problems ("storms")
- I wish to make it part of 7.x

### Conclusion

- No performance regression against previous model
- Lower interrupt latency
- Less interrupt problems ("storms")
- I wish to make it part of 7.x