# Linaro Connect Q4'12

# Peeking into your Linux System

Vincent Guittot <vincent.guittot@linaro.org>
Linaro Power Management Working Group

# Tools

- Used during this sessions
  - ftrace / trace-cmd
  - kernelshark
  - cylictest
  - sysbench
  - taskset
  - IRQ affinity
  - arm-probe

Vincent Guittot
Linaro Power Management Working Group

# Hardware

- TC2 board

  - Access to big and LITTLE power domain

- ARM probe

  - Up to 3 channels
  - 10kHz sampling rate

# Software

- Full Ubuntu Image 12.09

- Linaro ARM Landing Team kernel

- Disable HMP task placement
  - Disable load balance between cluster

- Additional patch set

Vincent Guittot
Linaro Power Management Working Group

# Goals

- Spy system behavior

- Understand wake up sources

- Understand scheduling behavior

- Exercise your system

# ftrace overview

- ftrace

  - In kernel tracing utility

  - Trace specific events like entering an idle state..

  - Trace function call

  - And more ...

- Location

  - <debugfs path>/tracing/

Vincent Guittot
Linaro Power Management Working Group

# ftrace overview

- Useful articles in addition to Documentation:
    - http://lwn.net/Articles/366796/
    - http://lwn.net/Articles/365835/


- Use trace-cmd tool
    - Simplify ftrace configuration and use

Vincent Guittot
Linaro Power Management Working Group

# trace-cmd overview

- trace-cmd encapsulates ftrace

  - Package available for Ubuntu image

  - Installed in next android image


- Graphic viewer : kernelshark

  - Nice human readable interface

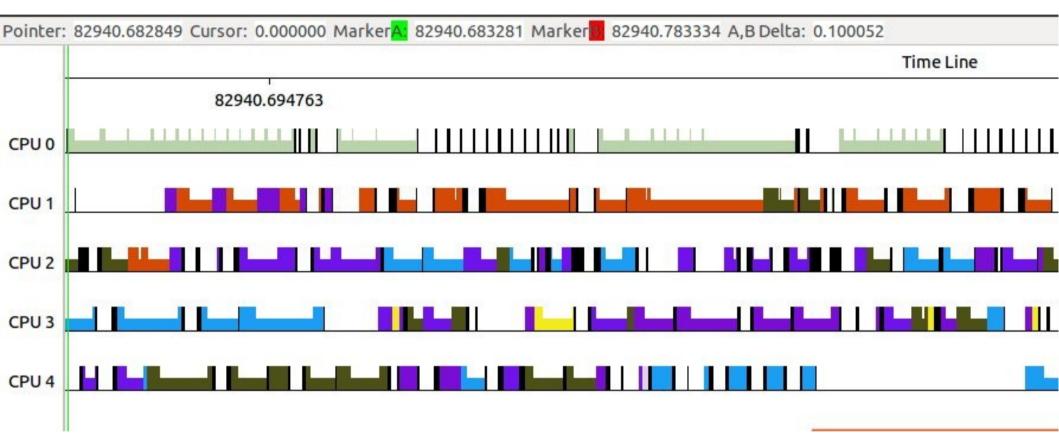  - Package available for Ubuntu image

# Start to trace

- Have a look at trace-cmd
  - Let make 1st traces

Vincent Guittot
Linaro Power Management Working Group

# trace-cmd record



## *trace-cmd record -e all -o example0.dat sleep 2*

- Huge activity...
- Tracing activity...

Vincent Guittot
Linaro Power Management Working Group

# •What to trace ?

- Core state:
  - Idle state
  - Frequency scaling

- Wake up source:
  - Interruption

- Activity:
  - Sched
  - Timer
  - Workqueue

Vincent Guittot
Linaro Power Management Working Group

# Ring buffer

- ftrace uses a ring buffer
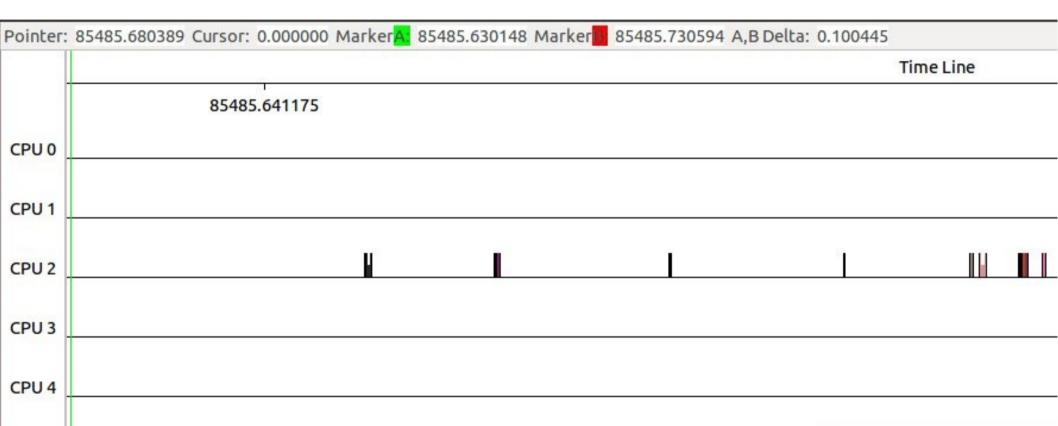
    - Limited number of events

- trace-cmd record

    - Periodically read the ring buffer

    - -s option set the period (default is 1ms)

    - Generate "spurious" activity : 1 process per core

    - Trig deferrable activity

# Sample period



***trace-cmd record -s 100000 -e irq -e sched -e timer -e workqueue  -o example1.dat sleep 2***

- Less activities ...

- Tracing activity ...

Vincent Guittot
Linaro Power Management Working Group

# Stay quiet

- trace-cmd start / stop / extract

  - No noise activity during the record

  - circular buffer

  - Increase the buffer's size
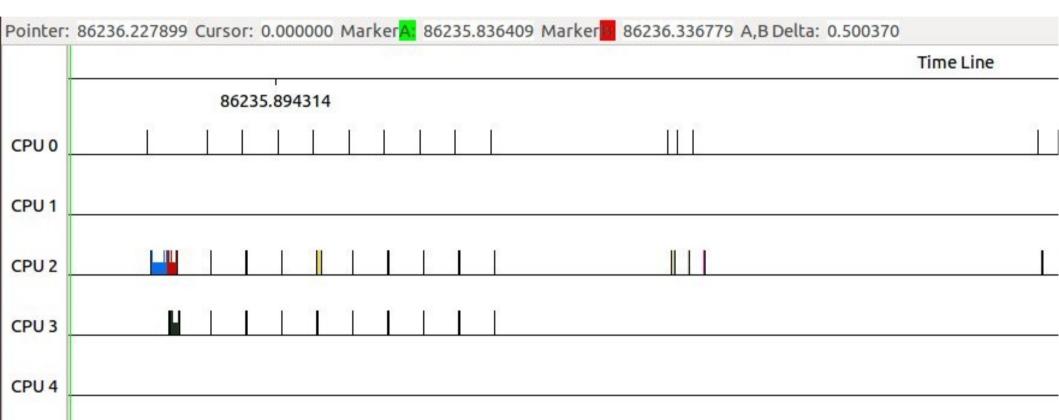

- Make a new trace

Vincent Guittot
Linaro Power Management Working Group

# Quiet trace



***trace-cmd start -b 4000 -e irq -e sched -e timer -e workqueue;
sleep 2; trace-cmd stop; trace-cmd extract  -o example2.dat***

- few activity...

- IRQ on big → CPU0

Vincent Guittot
Linaro Power Management Working Group

# IRQ affinity

- The real wake up source

- Influence where task will run

- Use 1st CPU of the mask

  - CPU0 → big core

- Change the affinity

  - /proc/irq/*/smp_affinity

- irqbalance daemon

- Set default affinity on LITTLE

# IRQ affinity



*trace-cmd start -b 4000 -e irq -e sched -e timer -e workqueue; sleep 2; trace-cmd stop; trace-cmd extract -o example3.dat*

- Nearly nothing on big cores...

# RT tasks

- On previous trace

  - Only some RT tasks on big

- There is not much we can do

  - Set task affinity with taskset

# Low CPU load

- Idle is almost aligned with what we want

- Test behavior with low load tasks

- Cyclictest
  - Create small tasks that wake up periodically

Vincent Guittot
Linaro Power Management Working Group

# cyclictest



*trace-cmd start -b 4000 -e irq -e sched -e timer -e workqueue
-e cpu_idle; cyclictest -t 7 -q -D 1; trace-cmd stop; trace-cmd
extract -o exampl4.dat*

- Stay on LITTLE cores

Vincent Guittot
Linaro Power Management Working Group

# Small task packing

- Most of cyclictest thread on CPU2

    - Some on CPU3

    - Few on CPU4
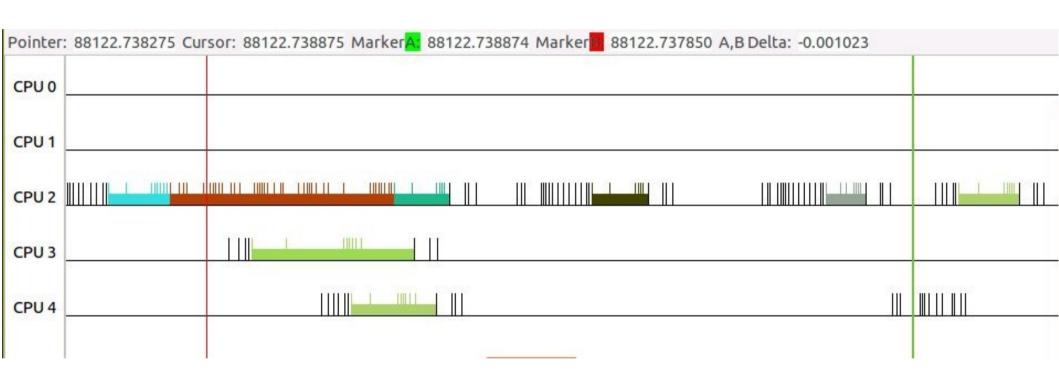
- Why ?

    - Pack buddy CPU is busy

    - Go back to default behavior

    - Worth to parallelize in a cluster

# Deeply in the trace



Pointer: 88122.738275 Cursor: 88122.738875 MarkerA: 88122.738874 MarkerB 88122.737850 A,B Delta: -0.001023

- Packing task → green line
- Spreading task → red line

Vincent Guittot
Linaro Power Management Working Group

# Power consumption

- Measure power of the use case
  - ARM-probe HW
  - Command line SW for acquisition
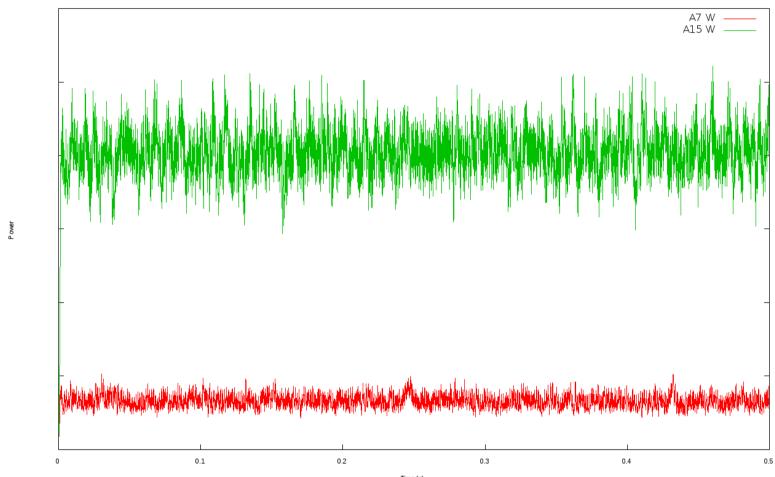    - Thanks to Andy Green
  - Gnuplot for displaying


- Make a measure

Vincent Guittot
Linaro Power Management Working Group

# 1st power result



*cyclictest -t 7 -q -D 1*

- Both cluster are always on and consuming !!!

Vincent Guittot
Linaro Power Management Working Group

# cyclictest

- Set the cpu_dma_latency QoS

  - Used by cpuidle

  - -e option set the QoS value

- Make another measure

Vincent Guittot
Linaro Power Management Working Group

# New power result



*cyclictest -t 7 -e 1000000 -q -D 1*

- Big cluster is off

- LITTLE cluster is on

Vincent Guittot
Linaro Power Management Working Group
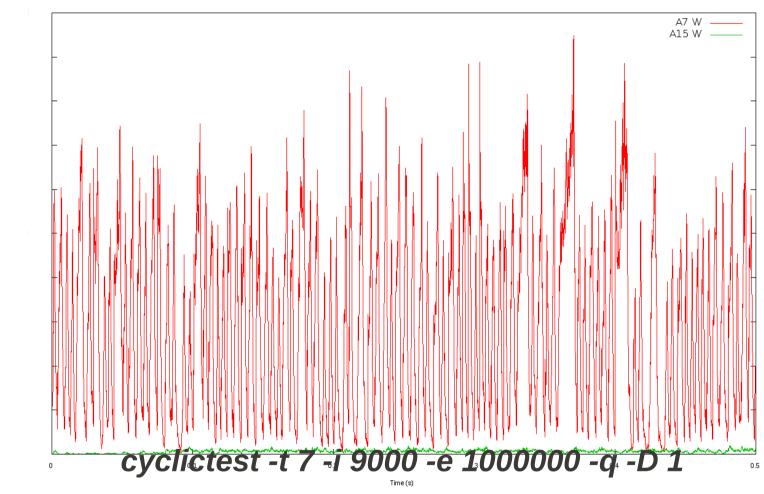
# cyclictest

- Cluster goes down iff idle duration >= 1ms

- Default cyclictest interval is 1ms

  - Can't enter off state

  - -i option increases the interval

- Make another measure

Vincent Guittot
Linaro Power Management Working Group

# Power measure



*cyclictest -t 7 -i 9000 -e 1000000 -q -D 1*

- Big cluster is off

- LITTLE cluster toggles

Vincent Guittot
Linaro Power Management Working Group

# Heavy task

- Small tasks are packed

- What about heavy task ?

- Sysbench
  - --test=cpu : prime number computation
  - --test=threads : lock thread

Vincent Guittot
Linaro Power Management Working Group

# sysbench



*trace-cmd start -b 4000 -e irq -e sched -e timer -e workqueue -e cpu_idle; sysbench –test=cpu –num-thread=1 –max-time=1 run; trace-cmd stop; trace-cmd extract -o exampl4.dat*
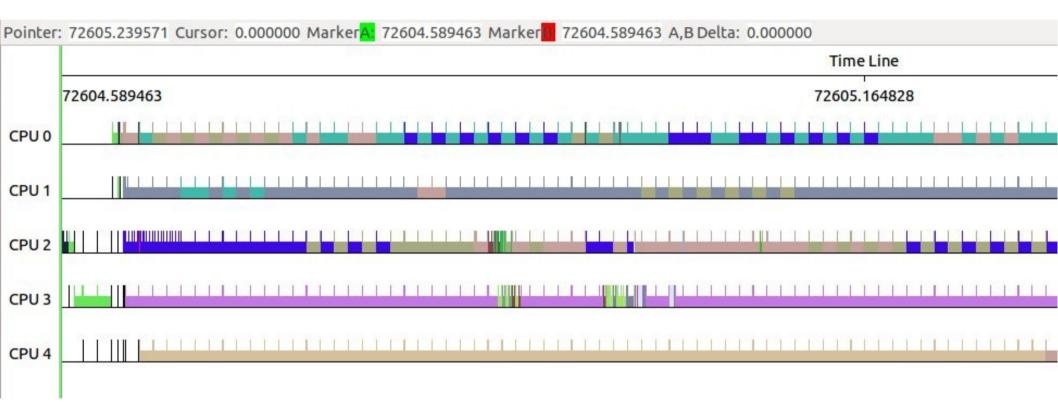
- Move on big core

Vincent Guittot
Linaro Power Management Working Group

# sysbench



***trace-cmd start -b 4000 -e irq -e sched -e timer -e workqueue -e
cpu_idle; sysbench –test=cpu –num-thread=7 –max-time=1 run;
trace-cmd stop; trace-cmd extract  -o exampl8.dat***

- Move on big core

Vincent Guittot
Linaro Power Management Working Group

# taskset

- Pin a task on a subset of CPU
  - Check what is the best scheduling behavior

- Performance
  - Fastest configuration

- Powersaving
  - Most thrifty configuration

# taskset

- Previous test

    - Better to put everything on big ?
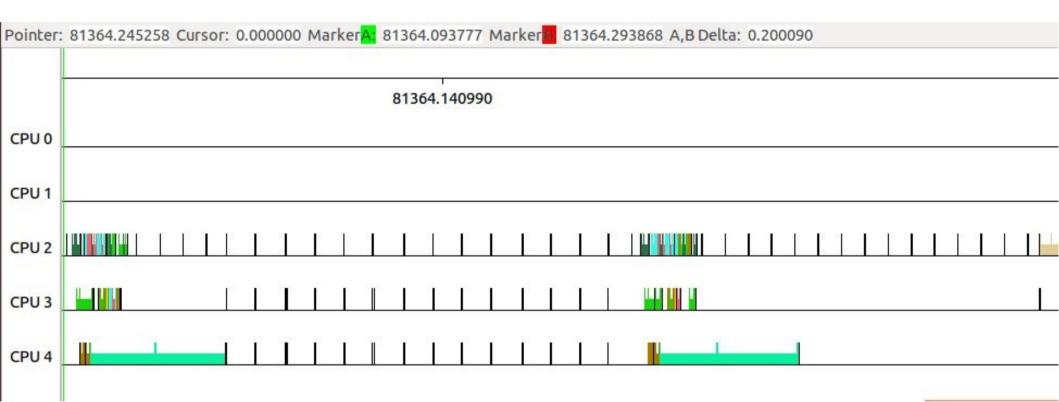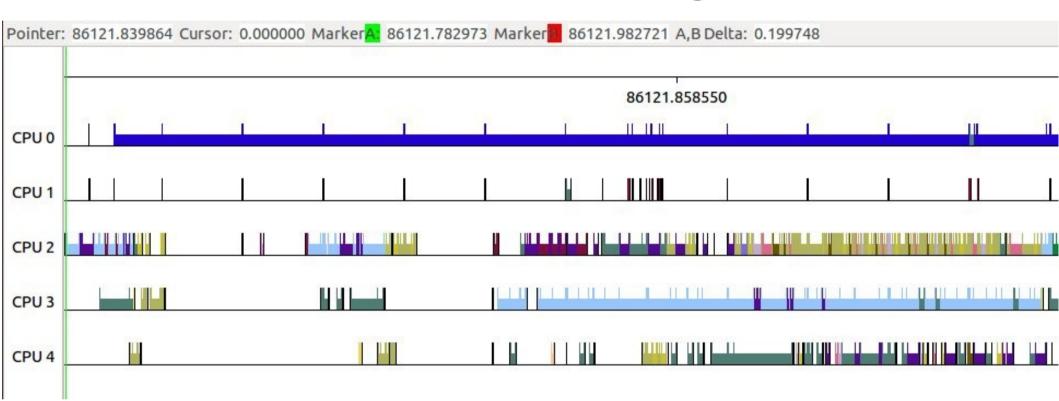
# Real use cases

- MP3 playback

- Web browsing

Vincent Guittot
Linaro Power Management Working Group

# MP3 playback

# Web browsing

Vincent Guittot
Linaro Power Management Working Group

# Question ?

Vincent Guittot
Linaro Power Management Working Group

# Thank you

Vincent Guittot
Linaro Power Management Working Group

# Backup slide

- MP3 sequence