

# LINUX Plumbers Conference

Richmond, Virginia | November 13-15, 2023



# Performance Monitor Control Unit (PMCU)

Jie Zhan - Huawei





# Contents

- 1. Motivation
- 2. What PMCU does?
- 3. Software Design
- 4. Short Demo
- 5. Main Challenge







# **Motivation**

• Performance Monitor Unit (PMU) are...

1. Counters for hardware events, e.g. cycles, cache misses, branch mispredictions 2. Great for performance analysis, resource management, debugging, task profiling, power modelling/management

BUT we found it...

## **1. Affecting latency/time-sensitive workloads**

Example: "perf stat" monitoring VMs, 100+ events, running rt-tests

Impact: cyclictest latency increases by 50%

### **2. CPU utilization increases with targets**

The more processes / threads we monitor, the more CPU utilization it occupies (0.1~0.2% CPU util / process)

<pre>root@localhost:/# while :; do cyclictest -D</pre>	20s -q;	sleep	) 10s;	done	2	0
<pre># /dev/cpu_dma_latency set to Ous</pre>					. A	no perf
T: 0 ( 694) P: 0 I:1000 C: 20000 Min:	11 Act:	12	Avg:	13	Max:	88
<pre># /dev/cpu_dma_latency set to 0us</pre>						
T: 0 ( 697) P: 0 I:1000 C: 20000 Min:	11 Act:	22	Avg:	12	Max:	26
<pre># /dev/cpu_dma_latency set to 0us</pre>						
T: 0 ( 700) P: 0 I:1000 C: 20000 Min:	11 Act:	12	Avg:	12	Max:	25
<pre># /dev/cpu_dma_latency set to 0us</pre>						
T: 0 ( 703) P: 0 I:1000 C: 20000 Min:	11 Act:	12	Avg:	12	Max:	24
<pre># /dev/cpu_dma_latency set to 0us</pre>						
T: 0 ( 706) P: 0 I:1000 C: 20000 Min:	6 Act:	12	Avg:	12	Max:	92
<pre># /dev/cpu_dma_latency set to 0us</pre>						
T: 0 ( 709) P: 0 I:1000 C: 20000 Min:	17 Act:	18	Avg:	18	Max:	36
<pre># /dev/cpu_dma_latency set to 0us</pre>						
T: 0 ( 712) P: 0 I:1000 C: 20000 Min:	17 Act:	21	Avg:	18	Max:	91
<pre># /dev/cpu_dma_latency set to 0us</pre>						
^CT: 0 ( 715) P: 0 I:1000 C: 1309 Min:	11 Act	t: 1	.2 Avg	: 1	.2 Max:	perf 28
^C						

### Latency impact on cyclictest by PMU monitoring

%CPU	%MEM	TIME+	COMMAND
0.2	Θ.Θ	0:16.02	perf
0.2	Θ.Θ	0:15.78	perf
0.2	0.0	0:14.69	perf
0.2	0.0	0:02.33	perf
0.2	0.0	0:02.21	perf
0.2	0.0	0:02.16	perf
0.2	0.0	0:02.16	perf
0.2	0.0	0:02.13	perf
0.2	0.0	0:02.13	perf
0.2	0.0	0:02.05	perf
0.2	0.0	0:01.99	perf
0.2	0.0	0:02.00	perf
0.2	0.0	0:01.93	perf
0.2	0.0	0:01.94	perf
0.1	0.0	0:02.43	perf
0.1	0.0	0:02.26	perf
0.1	0.0	0:02.24	perf
0.1	0.0	0:02.05	perf
0.1	0.0	0:02.06	perf
0.1	0.0	0:02.09	perf

CPU utilization occupied by "perf"







# What PMCU does?

A module that controls core PMUs through external memory-mapped interfaces, offloading monitor work from CPU

### Function

Configure core PMUs, switch events, and save PMU readings to memory

## Procedure

- Process an event list for a few rounds
- In each round:
  - Read a set of event IDs and set to PMUs for all cores
  - Start counting, wait for a period, stop counting
    Save all PMU readings to memory
- Issue an IRQ when all rounds are done

### Hardware view of PMCU (simplified)







# **Software Design**

Designed with perf\_event auxtrace framework

4 layers:

- sysfs: configs, e.g. pass event IDs
- perf tool perf-record controls PMCU data recording perf-report/script decodes PMCU data
- perf\_event framework allocates AUX buffers, mapped to user space, passed to driver
- driver

interacts with the framework, configures hardware to produce results into the AUX buffers

Reference: <u>https://lwn.net/Articles/922351/</u>

\_inux Plumbers Conference | Richmond, VA | Nov. 13-15, 2023

# **Short Demo**

1. Enter event numbers in the 'user\_events' file:

echo "11 8 b 3c 24 23" > /sys/devices/hisi\_pmcu\_sccl3/user\_events

2. Start the sampling with 'perf-record':

perf record -e hisi\_pmcu\_sccl3/nr\_sample=10,sample\_period\_ms=4/

3. Decode PMU data with 'perf-script':

## perf report -D



**Organization of PMCU data output** 

### **Organization of PMCU data output**

Header

HISI PMCU data: size 0x7840 bytes												
Header: s	size 0x4	0 byte	s									
00000000: 0	00 00 40	00 06	00 00	00	30	00	00	00	00	00	00	00
00000010: 0	00 Oc 00	00 01	00 00	00	06	00	00	00	11	00	00	00
00000020: 0	08 00 00	00 Ob	00 00	00	3c	00	00	00	24	00	00	00
00000030: 2	23 00 00	00 00	00 00	00	00	00	00	00	00	00	00	00
Auxtrace but	ffer max	size:	0x400	0000								
Number of PM	MU count	ers in	para	lel	: 6							
Number of mo	onitored	CPUs:	48									
Compatible r	mode: no											
Subsample size: 0xc00												
Number of subsamples per sample: 1												
Number of ev	vents: 6											
Event 0: 0	0x0011											
Event 1: 0	8000x0											
Event 2: 0	0x000b											
Event 3: 0	0x003c											
Event 4: 0	0x0024											
Event 5: 0	0x0023											
Data: siz	ze 0x780	0 byte	s									

### Data body

. . .

00000300:	00000000009190c8	Event	0011	CPU	38
00000308:	0000000000905798	Event	0011	CPU	46
00000310:	0000000000009905	Event	0011	CPU	7
00000318:	0000000000917760	Event	0011	CPU	15
00000320:	0000000000000a052	Event	0011	CPU	23
00000328:	0000000000000000	Event	0011	CPU	31
00000330:	000000000000c8b7	Event	0011	CPU	39
00000338:	00000000000c5d14	Event	0011	CPU	47
00000340:	000000000015ceb9	Event	0008	CPU	0
00000348:	000000000015cd5f	Event	0008	CPU	8
00000350:	00000000018e7a68	Event	0008	CPU	16
00000358:	0000000000007394	Event	0008	CPU	24
00000360:	00000000001e51ae	Event	0008	CPU	32
00000368:	0000000000e7212c	Event	0008	CPU	40
00000370:	00000000017468fb	Event	0008	CPU	1
00000378:	0000000001573094	Event	0008	CPU	9
00000380:	000000000009f40	Event	0008	CPU	17
00000388:	00000000180390c	Event	0008	CPU	25







# Main Challenge

# **1. Synchronization between external and internal PMU accesses**

ARM PMUs can be accessed from both internal sysreg (from local CPU) and external memory-mapped interfaces (from other CPUs or devices, e.g. PMCU) simultaneously. However, a mechanism for synchronizing counter accesses or reserving counters is missing.

### **2. Kernel support for ARM PMU external memorymapped accesses?**

ARM PMUs are currently controlled from internal sysreg in Linux. External memory-mapped interfaces are not yet enabled.

Shall we enable support for it? Prioritized cores could avoid handling perf events and focus on the main workload.



# Plumbers Conference

Richmond, Virginia | November 13-15, 2023

