

Performance Oriented Firmware Implementation of Error Handling on OCP Platforms for Large Scale Data center

Zhengyu Yang, Hardware System Engineer, Facebook Chris Davis, Core Data Production Engineer, Facebook





Hardware Management









- MANAGEMENT
- Generic performance problem on OCP platforms
- Problem diagnose and generic
- firmware solution
- Results and verification strategy
- Summary



White Papers









Performance Issues Affects Various OCP Platforms Platforms Affected in Facebook Fleet







Performance Issue Degrades Availability & Latency Symptoms Observed in FB production

- Low latency systems like caches experience timeout spikes
- Failures in time critical processes like DB master promotion
- Packet loss, especially impacting new connection attempts
- Higher latencies & error rates in general



имит





Diagnosing & Detecting System Stalls System Unresponsiveness (a.k.a System Stall)

- Machine functions correctly but stalls for hundreds or even thousands of milliseconds at a time
- "stall" means the machine does 0 work for this time period

										_									
Task		Runtime r	ns	l Sw	itches	I	Average	delay	ms	I	Maximur	ı delay	ms	Mo	aximum	delay	at		I
ksoftirqd/0:3	I	261.106	ms		166	I	avg:	25.445	ms	I	max:	570.087	ms	m	ıx at:	12725	714.7	'54445	s
<pre>rcu_sched:8</pre>		640.922	ms		8048		avg:	0.108	ms	I	max:	381.426	ms	mo	ix at:	12725	671.2	19097	' s
<pre>scribe_cat:(99)</pre>		2739.591	ms		1171	I	avg:	0.349	ms	I	max:	269.621	ms	mo	ix at:	12725	714.4	54340	S
xargs:(18)		2541.290	ms		22096	I	avg:	0.018	ms	I	max:	254.727	ms	mo	ix at:	12725	683.0	15549	S
kworker/3:0:244396		2.960	ms		121	I	avg:	2.193	ms	I	max:	254.498	ms	mo	ix at:	12725	691.3	76932	s
ksoftirqd/3:29		3.871	ms		99	I	avg:	2.751	ms	I	max:	254.495	ms	mo	ix at:	12725	691.3	876942	s
cat:(21992)		23824.856	ms		66206	I	avg:	0.020	ms	I	max:	254.491	ms	mo	ix at:	12725	690.0	99208	s
carbon-global-s:(3)		446.988	ms		492	I	avg:	1.647	ms	I	max:	254.489	ms	mo	ix at:	12725	691.3	76856	S
LadDxChannel:2861439		660.121	ms		5147		avg:	0.065	ms		max:	254.454	ms	mo	ix at:	12725	691.3	376797	' s
<pre>FutureTimekeepr:(146)</pre>		1963.384	ms		6555	I	avg:	0.283	ms	I	max:	254.442	ms	mo	ix at:	12725	691.3	76753	s
<pre>[xarexec] /usr/:(21)</pre>		14159.697	ms		9236	I	avg:	0.036	ms	I	max:	254.427	ms	mo	ix at:	12725	691.3	876745	s
kworker/1:0:1303233		2.927	ms		131	I	avg:	2.058	ms	I	max:	254.426	ms	m	x at:	12725	691.3	76798	s
ksoftirqd/1:17		5.037	ms		109	I	avg:	2.447	ms		max:	254.412	ms	mo	ix at:	12725	691.3	876817	s





Diagnosing & Detecting System Stalls Detection

- Detection process simply sleeps and measures how long it actually slept.
- Stalls manifest as a significant delay in waking up.
- Doing this on all cores increases signal





Script for Detecting Stalls

```
while True:
 start = time()
 sleep(100ms)
 elapsed = time() - start
 if elapsed > 150ms:
    print("Stall Detected")
```



System Stalls Correlate with Hardware Errors Correctable Errors are common

Error]:	{1861667}[Hardware	[12725978.954807]
Error]:	<pre>{1861667}[Hardware</pre>	[12725978.972915]
Error]:	<pre>{1861667}[Hardware</pre>	[12725978.991358]
Error]:	<pre>{1861667}[Hardware</pre>	[12725979.003874]
Error]:	<pre>{1861667}[Hardware</pre>	[12725979.016388]
Error]:	{1861667}[Hardware	[12725979.029095]
Error]:	{1861667}[Hardware	[12725979.041610]
Error]:	{1861667}[Hardware	[12725979.052407]
Error]:	{1861667}[Hardware	[12725979.066313]
Error]:	{1861667}[Hardware	[12725979.078838]
Error]:	{1861667}[Hardware	[12725979.088590]
Error]:	<pre>{1861667}[Hardware</pre>	[12725979.100438]
Error]:	<pre>{1861667}[Hardware</pre>	[12725979.115221]
Error]:	<pre>{1861667}[Hardware</pre>	[12725979.126871]
	Error]: Error]: Error]: Error]: Error]: Error]: Error]: Error]: Error]: Error]:	<pre>{1861667}[Hardware Error]: {1861667}[Hardware Error]: </pre>



JMMIT

ware error from APEI Generic Hardware Error Source: 0 nas been corrected by h/w and requires no further action nt severity: corrected ror 0, type: corrected ection_type: PCIe error ort_type: 4, root port ersion: 1.16 ommand: 0x0010, status: 0x0547 evice_id: 0000:00:01.0 Lot: 0 econdary_bus: 0x01 endor_id: 0x8086, device_id: 0x6f02 ass_code: 040600 ridge: secondary_status: 0x2000, control: 0x0013



Breakdown of HW Error's Performance Impacting Aspects









Hardware Error Handling Overview Stack view

- Components
 - Detect, Report and Remediate
- Contributors/Owners
 - **OS & Software** : Harvest failure information and apply remediation method if any
 - Firmware : Implement HW error handlers
 - Hardware: Provide HW capabilities and error information for error handling









SMM and SEL

- System Management Mode (SMM)
 - Contains SMM HW error handlers
 - Context switch all cores when entered
 - Suppose to be fast
- System Event Log (SEL)
 - Industry standards, allow content customization
 - Created and sent out by BIOS using IPMI commands within SMM
 - IPMI commands can be expensive









Performance Issue Root Cause System Stall due to SMM & SEL



JMMIT



all	Stall	

Generic problem due to FW, across OCP platforms!

Long latency for single PCIe error logging

Performance Focused Hardware Error Handling FW Solution Strategy

- Failure event based reporting vs. Error occurrence based reporting
- Consolidated, Standardized, Scalable error record
- Overhead minimized error information communication









Performance Focused Hardware Error Handling FW Solution Example -- BIOS

- Failure Event
 - High Frequency Error Burst
 - Low Frequency Error Tracking
- Unified failure record format







Performance Focused Hardware Error Handling FW Solution Example – Bridge IC



High-level Flow of HW Error SEL Logging with BIC





Performance Focused Hardware Error Handling FW Solution Example – BMC

- Overhead minimized error information communication
 - Ex. Enhanced BMC SEL handling flow -- "Ack Quick"*



* Submitted to oBMC committee and pending upcoming spec adoption



Run in Background







Performance Focused Hardware Error Handling FW Solution Example Results from 3 OCP platforms









FW Implementation Verification Strategy

- Error Injection
 - HW capability support
 - Error type coverage
- Performance focused
 - Use case driven metrics
 - Quantifiable criteria
- Verification model
 - Performance measured with hardware failures/errors
 - Error occurrences' randomness



лими





Summary

• Hardware error handling caused system performance impact is a common problem across OCP platforms

common firmware layer and applicable to various OCP platforms

needs to be system performance aware



имит

Improvements (ex. logging strategy, processing efficiency) can be made at the

Development and validation of firmware solutions for hardware error handling



Call to Action

- Contribution to OCP (Planned Q3 2019)
 - IPMI SEL HW Error record format and content
 - BIC FW "Return First" Model
 - BMC FW "Ack Quick" SEL processing flow
- Check the fleet

 - OCP system stall in the fleet • Correlate system stall with hardware error occurrences Adopt the presented fixing strategy to see
 - improvements



UMMIT







Open. Together. OCP Global Summit | March 14–15, 2019





