

13th ANNUAL WORKSHOP 2017

LDMS AND INFINIBAND @ SANDIA

Serge Polevitzky, SAIC HPC Support at SNL

SAIC

[March, 2017]



SAND Number: SAND2017-2862

INTRODUCTION

What You Should Expect

- Some Disclosures: System Described, Presupposition, etc.
- Next, Why We Started Our Investigation
- Background on Tools Used
- What We've Seen (The Recent Story)
- Expect "Support" from the Audience ! [Interrupt to Correct, to Add Info, *etc.* ... in short, PLEASE PARTICIPATE]
- We Will Follow the old Cray Research Mantra of "We Don't Take Ourselves Seriously, But We Take What We Do Very Seriously."



THE PROBLEM / THE "INSURMOUNTABLE OPPORTUNITY"

FABRIC MELTDOWN

One (More ?) IB-Connected Clusters Had Meltdowns

- Started w/ Experiencing an Occasional IB "Meltdown" ["Cure" Is/Was Pretty Draconian]
- We Knew of Livermore's "Sniper Script" Solution, But ...
- Didn't Have Sufficient Evidence as to Whom to "Snipe"
- Tribal Knowledge Was That It Was Our LNET Routers
- Frequency of Meltdowns Started to Increase
- A "Must Resolve This <u>Now</u>" Edict Was Issued

- Relocated & Reduced Number of LNET Gateways (...?) &
- Relocated the Subnet Manager (SM)
- An *a priori* Solution
- And It Worked ! Life Got a Lot Better. We Took Stock.
- Our "Tool Kit": Smart People (coup d'oeil) !, perfquery, PerfMgr, LDMS
 - perfquery (But Many Functions are optional !?)
 - PerfMgr, But the Data Wasn't "Helpful"
 - LDMS. Just Starting to Be Used, so Little "Street Cred"

~ 18 perfquery commands are "optional" (YMMV)

- -D, --xmtdisc show transmit discard details. This is an optional counter.
- -E, --rcverr show receive error details. This is an optional counter.
- --slrcvfecn show SL Rcv FECN counters. This is an optional counter.
- --slrcvbecn show SL Rcv BECN counters. This is an optional counter.
- perfquery –rcvcc perfquery: iberror: failed: cannot query PortRcvConCtrl

- We "Benched" perfquery. This Left PerfMgr & LDMS.
- Three "Sticking Points" Emerged w/ PerfMgr & LDMS.
 - 1) PerfMgr found 1 Link_Downed; LDMS Saw Many,
 - 2) PerfMgr Saw Many VL15 Drops; LDMS Saw 0,
 - 3) PerfMgr Saw 0 port_transmit_waits; LDMS Saw Many.
- port_transmit_wait :: Great Indicator of Congestion ! The Canary in the Coal Mine (for congestion).
- Let's Start w/ 1) LDMS Sees Many Link_Downed Counts

- Link_Downed Counts: PerfMgr & LDMS ...
 - Both Saw the Link_Downed Happening at the Same Time (±)
 - LDMS Sees the Link_Downed @ 08:25 (505 minutes after 12:00)
 - Why is 505 Minutes After Midnight Important ? [See Next Chart]
 > count_ib.link_downed <- sum(data_set\$ib.link_downed)
 - > count_ib.link_downed
 - > [1] 935

scatter chart from LDMS: sysclassib.1487055600 [Feb 14 2017], starting at 00:00:00 thru Feb 14 2017, 23:59:00 MT (0 dropouts)



Platform:x86_64-apple-darwin13.4.0; Host R version 3.3.2 (2016-10-31); Nickname:Sincere Pumpkin Patch; Script running:/Users/ LDMS_xmit_wait_all_scatter_44n1.sh; Script executed:February-16-2017 13:37; Executed by:

- Invoking the Willing Suspension of Disbelief ...
 - PerfMgr Saw 1 Link_Downed [the Court accepts this Claim]
 - Straight from LDMS .csv file for Link_Downed Event \rightarrow
 - ✓ Time … CompID … ib.link_downed
 - ✓ 1487085900**120007
 - REDACTED-HOST-NAME-login[1-8] ldmsd_idbase=120000 [login node 7 Had the 1 Link_Downed Count]
 - Submit to the Court that PerfMgr & LDMS Agree; LDMS PerfMgr Just Reporting the Data Differently

Sticking Point #1 Resolved ...

- Continuing the with the Willing Suspension of Disbelief:
- Sticking Point #2 (PerfMgr Sees Many VL15s Dropped ...) and LDMS Sees ZERO
- PerfMgr Looks at Switches and HCAs ! LDMS Only Looks at HCAs (at Least at the Time of These Data Captures)
 - 1) How Can an HCA Know It has Dropped a VL15?
 - 2) Conclusion: Since PerfMgr Sees Switch Info, "Naturally" PerfMgr Will See the VL15s Dropped & LDMS Will Not
- Submit that Sticking Point #2 Is Resolved

- More Willing Suspension of Disbelief Now Required:
- PerfMgr Sees ZERO port_transmit_waits (from Any Source)
 - LDMS Shows → Next Slide [a One-Hour Capture"]
 - ✓ Whatever It Is, It Gets Reset Every 10-Minutes
 - Some (Possibly Small) Increment Every Minute
 - The Counts Are Non-Zero (so Disagreement w/ PerfMgr)

Next Slide → {Many port_transmit_waits} // PerfMgr Sees 0







Platform:x86_64-apple-darwin13.4.0; Host Rversion 3.3.2 (2016-10-31); Nickname:Sincere Pumpkin Patch; Script running/Users Desktop/ Lob_14_2017/ Lob_14_2007/ Lob

OpenFabrics Alliance Workshop 2017

- PerfMgr Has a Default On/Off Switch & A Threshold Switch Too
 - 1) PerfMgr Has a Primary Switch that says "Don't Report Anything"
 - 2) "perfmgr_xmit_wait_log", OPT_OFFSET(perfmgr_xmit_wait_log), opts_parse_boolean, NULL, 0 }
 - 3) perfmgr_xmit_wait_log wasn't defined in our opensm.conf
 - 4) So PerfMgr (for us) Was Blind to All port_transmit_waits
 - 5) And a 2° Switch, a Threshold for "Do Not Report Count if Less Than X" [#FFFF]
 - 6) We Then Turned the Switch to "On," and Set the Threshold to 0
 - And

PerfMgr (cont'd)

- 7) No Change ! Still Seeing 0 port_transmit_waits from PerfMgr ...
- 8) Still Seeing large Number of port_transmit_waits from LDMS
- 9) SME #1: port_transmit_wait Counts Turned Off Because the Counts Were so Large & It's a 32-bit Counter ...
- 10) SME #2: "The scatter plot looks reasonable for port_transmit_ wait counts"
- 11) Hmmm ... ??
- 12) For the Moment, Suggest Leaning Toward Assuming the LDMS Data Has Some Merit

Sticking Point #3 Still NOT QUITE Resolved.

- Still Don't Know What Action to Take for Given Any Given port_transmit_wait Counts, Even If We Were to Have Total Faith in Those Counts from Either PerfMgr and/or LDMS.
- Again, port_transmit_wait Is the Canary in the Coalmine for Congestion, so Monitoring It Would Seem to Be Very Important
- "In God We Trust. All Others Must Bring [Relevant] Data."
- This IB Exploration Is Very Much a Work in Progress. But Here Is a Brief "exploration" of LDMS Data from the One Sandia HPC Cluster



LOOKING AT THE DATA (GIVEN WHAT WE KNOW)

IN AN IMPERFECT WORLD ...

Suggested Strategy

• First, See Any IB Errors ? If so, Attend to Them.

If No Errors, Look at Performance, "Got Congestion ?"

 Performance & Congestion May Take Some Delving into, Not as Clear-Cut a Process as with Error Checking



Workshop 2017



Platform:x86_64-apple-darwin13.4.0; Host Rversion 3.3.2 (2016-10-31); Nickname:Sincere Pumpkin Patch; Script running:/Users Desktop/ eb_26_201







Total port_transmit_wait counts for all 1263 nodes, all nodes plotted by the minute

Maximum port_transmit_wait count at 01:20:00 sktop/_____feb_26_2017/____LDMS_xmit_wait_all_scatter_46g.sh; Script executed:February-28-2017 10:40; Executed by: Time in Hours ----> Platform:x86_64-apple-darwin13.4.0; Host: R version 3.3.2 (2016–10–31); Nickname:Sincere Pumpkin Patch; Script running:/Users

Vorkshop 2017

File processed:sysclassib.1488092400

CONCLUSION - 1

- Our Understanding of IB Is Still "New" ...
 - Spread SM, LDMS, & LNET GWs Across Switches, Nodes
 - If Gateway Congestion, Increase Receive Buffers (4X ?)
 - Increase the Number of VLs ... [Decreases Each VL's Buffer]
 - More LNET Gateways Probably Better than Fewer
 - Use Both PerfMgr & LDMS
 - Use SPLUNK [Our PerfMgr & LDMS Feed into SPLUNK]

CONCLUSION - 2

Some "Possibles" & "Requests"

- Have a Metric with Known Performance Characteristics (port_transmit_wait Numbers); On Bring Up, Run a Quick Job that Generates PerfMgr & LDMS Current Values
- Request a Rosetta Stone [*e.g.*, Clearly Define All Variables [What Is/Are "ib.COUNTER_SELECT2_F", "port_xmit_constraint_errors," *et al.* ?]
- Request Definitions of "Hidden Info" e.g., "Alert ! the port_transmit_wait Switch Is Set to NOT Report Data";
 "Alert! port_transmit_wait Values > X Usually Result in a port_transmit_discard," etc.



THANK YOU ! ... ANY OUESTIONS / DISCUSSION ?

Serge Polevitzky SAIC @ SNL







13th ANNUAL WORKSHOP 2017

BACKUP INFORMATION

Serge Polevitzky SAIC @ SNL





EXCEL.CSV INFO

Excel (on a Mac) Fails to Read in 320MB "Whole Day" — Can Read an Hour's Worth

Entry #	Line #	Time	Time Delta	Compld	Port_Xmit_Wait Count	Delta	
1	4	1488092400		130024	2,082,299		
2	1266	1488092460	60	130024	792,893	1,289,406	reset prior
3	2531	1488092520	60	130024	800,462	7,569	
4	3790	1488092580	60	130024	1,141,028	340,566	
5	5 5054	1488092640	60	130024	1,174,117	33,089	
6	6317	1488092700	60	130024	1,535,149	361,032	
7	7579	1488092760	60	130024	1,544,015	8,866	
8	8 8838	1488092820	60	130024	1,982,079	438,064	
9	10172	1488092880	60	130024	1,990,948	8,869	
10) 11434	1488092940	60	130024	2,147,879	156,931	
11	12654	1488093000	60	130024	2,580,935	433,056	
12	13959	1488093060	60	130024	344,518	2,236,417	reset prior

HANDY TOOL

Created: Thu Aug 18 18:28:04 UTC 2016, by processing file Erik_data_clean.txt



INTERESTING

scatter chart from LDMS: sysclassib.1489511440 [Mar 14 2017], starting at 11:12:00 thru Mar 14 2017, 23:59:00 MT (0 minute dropouts)



Total port_transmit_wait counts for all 1267 nodes, all nodes plotted by the minute

OpenFabrics Alliance Workshop 2017

Time in Hours ----> Maximum port transmit wait count at 23:59:00

INTERESTING

scatter chart from LDMS: sysclassib.1489511440 [Mar 14 2017], starting at 11:12:00 thru Mar 14 2017 , 12:13:00 MT (0 minute dropouts)



Time in Hours ----> Maximum port_transmit_wait count at 12:13:00

abrics Alliance Workshop 2017