Lustre Cluster Hera Prague - spring 2012

HEPiX

Hera – the IB lustre part

- IB based, commodity hardware (cheap :-))
 Even boot is over IB (low latency booting :-))
- Installation & configuration in a couple of days
 - thanks to our Chef based fabric management :-)
- Independent Ethernet network for IPMI Modules
- Connected via LNET router cluster to IP lustre
- "Pilot" Version with 2.400 disks
- Two independent "thermal measures"
 - System temperature based on IPMI Modul
 - => graceful shutdown if threshold reached
 - Rack temperature inside/outside door above thresh.
 - => emergency shutdown of PDU





Lustre at GSI: about 7000 disks, 6.5 PB Raw 150 Fileserver, about 1 Tb/s I/O



Prometheus / Hera - IB cluster in the Minicube

Executive summary: works! Success! :-)

• Prometheus: 10,000 core IB cluster pilot

 Hera:IB pilot with 2,400 disks lustre, 50 File servers, will be extended soon

Prometheus + Hera – first results

- Few weeks in operation now
- 1.000.000 physics jobs already finished
- 6 PB data throughput per week (few users ...)
- ~.5 Tbit/s I/O troughput already reached (in the Hera part, limited by the switched backplane)
- => Already (successfully) used by Alice, Acclerator-Physics, GSI-Theory, HIM

Prometheus + Hera – first results (out of the box :-))

- IB works reliable
- Speed measured with Prometheus
 - RDMA per node: 3.1 Gbyte/s (expected: Mellanox: ~3.1GB)
 - Parallel access to lustre per node: Read ~2 Gbyte/s about a factor 20 faster compared to our IP clients.....
- Fileserver performance now limited by the switched backplane..



HPC Test 450 Gbit/s (close to the backplane limit)



<u>File E</u> dit <u>V</u> iew Hi <u>s</u> tory <u>B</u> ookmarks <u>T</u> ools <u>H</u> elp														
🖕 🛶 🔻 🔁 🚳 🏠 💽 http://lxrm02.gsi.de:5000/report?display=gridengine%2Faccounting 🛛 🖓 🖬 Google														
🛅 Most Visit	ömlost Visited▼ 🥮 Getting Started 🔝 Latest Headlines▼													
Report			4											
Graph	Data F	Report Help												
Select repo	rts to display	: •												
			Cha	Canaal										
	,		Sil		,		,							
User	Jobs	CPU Time	Run Time	Memory Used	I/O	Jobs Max Resident Memory			Jobs Run Time					
				Over Time		<1GB	<2GB	<4GB	>4GB	<1h	<4h	<12h	>12h	
wiechula	11563	14kh	15kh	65.53PB/s	190.62TB	7521	3747	284	11	6317	4185	1061	0	
tneff	6	1h42m	5m37s	1.22TB/s	31.88MB	6	0	0	0	6	0	0	0	
vpenso	9	Os	8s	0.0B/s	6.27MB	9	0	0	0	9	0	0	0	
root	2	Os	2s	0.0B/s	1.21MB	2	0	0	0	2	0	0	0	
sma	192949	37kh	115kh	80.08EB/s	2.03PB	140600	51555	193	601	178022	14920	7	0	
jthaeder	104105	233kh	242kh	17.27EB/s	3.73PB	80305	0	23800	0	53333	26902	23869	1	
mfasel	20	3h58m	4h11m	27.88TB/s			Promet	neus Me	mory U	tilizati	on		R	
/report/gride	ngine/accoun	ting?format=htm	ıl		†				- 				TOOL	
					20 T									
										· / · · · · · ·	$ \left - \right $	FOBI		
$1 \in \mathbb{N}$ uays, $\mathbb{I} \in \mathbb{N}$ use 15.					10 т 🕌					به الم			Ē	
6 F	PB I/0	C		œ)				·····		<u>.</u>		IKE	
					<u> </u>					\ <u></u>		n.	20	
		0 +	Sun 00:00	Su	un 12:00	Мо	n 00:00	Mon 1	12:00					
	📕 Used Min 937.3GB Avg 2.1TB Max 7.3TB Last 2.2TB													
					Free	Free Min 1.6TB Avg 9.2TB Max 16.9TB Last 16.8TB								
Dono			L Cache	min I.	TID	AVG 8.5	10	"dX 12.41	D Las	st 1,2	в —			

Done

Next Steps: add more discs soon

110