IBM TechU

AIX Performance Tuning Top Tips

<text><text><image>







Patch order

- 1. Read the readmes
- 2. Usually do HMC first
- 3. Then server firmware
- 4. NIM server (it should be standalone and needs to be at highest level) plus I/O firmware
- 5. Then VIO servers and any I/O firmware
- 6. LPARs AIX, IBM i, Linux
- 7. Note I document everything
- 8. I write up every step of the upgrade before I do it and then tweak as I go along











Care and Feeding of VIOS - Jaqui Lynch

Efixes and ifixes

Many security patches are put on using efixes or ifixes

AIX and VIO servers need these to be applied – **use FLRTVC to determine what fixes are needed** Run flrtvc and download and install the ifixes that are needed https://www-304.ibm.com/webapp/set2/sas/f/flrt/flrtvc.html

You will need to read the .asc file to determine which efixes you need for the level you are running Efixes & ifixes are downloaded using FTP or http. Links are provided by FLRTVC

/usr/sbin/emgr –I lists them emgr –P lists the patches and the packages they affect To apply a fix change into the directory it is in and then run it in preview mode: cd /usr/local/soft/fixes/bind_fix17 emgr -p -e IJ25927s2a.200708.epkg.Z Remove the –p and run again for real if the preview was successful: emgr -e IJ25927s2a.200708.epkg.Z

If you run emgr –l and there are no fixes listed then you most likely have security holes that need patching, specifically Java, openssh and openssl. $$_{\odot}$$

9





3. Keep it simple and consistent

11

General

- Keep it simple
- Ensure LMB is the same on all servers if you want to use LPM
- Use hot pluggable adapters rather than built in ones Easier maintenance
- All adapters should be desired, not required
- Don't mix multipath drivers on HBAs
- Run HMC Scanner and/or Sysplan before and after all changes
- Plan for at least one update per year (IBM normally puts out 2)

General

- Mirror rootvg and page space if on internal disk
- Have a spare disk in your LPAR to use for cloning prior to updates
- NOTE AIX requires at LEAST 30GB in rootvg give it 100GB to 150GB
- Keep your rootvg very clean
- If you need user or application filesystems put them outside of rootvg
- Check page spaces
 - Make sure all page spaces are the same size and on different LUNs
- Add logging and set up dump devices properly
- Check errpt regularly
- NEVER run at 100% entitlement ensure it is high enough and there are plenty of VPs and memory
- Backup regularly use NIM or scripts
 - Backups should include mksysb images that can be used for bare metal restores, not just application data
 - Mksysbs should be taken for all AIX LPARs and VIO servers

13



Documentation is Critical

- Use HMCScanner and Sysplan
- Put together spreadsheet of documentation
- All equipment and serial numbers
- UAK expiration dates on servers (and AIX in the future)
- Customer numbers
- Server and I/O firmware levels, VIO and O/S levels
- HMC information including version, BMC and PNOR, networking etc
- IP addresses
- Resource profiles
- Adapter allocations
- Standards used for network, vSCSI, NPIV mappings
- Actual vSCSI assignments
- Actual NPIV vfcmaps
- Vfchosts and their associated WWPNs
- SEA and virtual ethernet VLAN assignments
- Switchports for SAN and network
- Power needs and PDU mapping
- Anything else you can think of

15



Use alt_disk_copy prior to changes	
 Have two disks so you can take a clone If rootvg is mirrored you will need to unmirror for maintenance No need to mirror if on SAN, always mirror rootvg if on internal disk 	
# lspv grep root hdisk0 00ce48c008314b9f rootvg active hdisk1 00ce48c03c8f2115 altinst_rootvg	
# bootinfo -b hdisk0	
exportvg altinst_rootvg alt_disk_copy –V –B –d hdisk1	
I always do a bosboot and rewrite the bootlist before any reboot Recovery if issues with the upgrade is to point the bootlist to the new	disk and reboot
If you want to use FBO (file backed optical) add a 3 rd disk in its o huge	own VG so rootvg does not get
Keep rootvg small and clean	
Don't forget mksysb images	17



What are you measuring?

- *Response time* is the elapsed time between when a request is submitted and when the response from that request is returned.
 - Amount of time for a database query
 - · Amount of time it takes to echo characters to the terminal
 - · Amount of time it takes to access a Web page
 - How much time does my user wait?
- *Throughput* is a measure of the amount of work that can be accomplished over some unit of time.
 - Database transactions per minute
 - File transfer speed in KBs per second
 - File Read or Write KBs per second
 - · Web server hits per minute
- Make sure you know the difference and what is important to you

19

What makes it go slow? - 1/2 • Obvious:- Not enough CPU • Not enough memory Not enough disk bandwidth Number of adapters · Queue depth and adapter queues · Number of disks Not enough network bandwidth Hardware errors - errpt Software errors - errpt Insufficient CPU entitlement Affects CPU performance · Affects network performance for SEA Too many or too few VPs Overallocating CPU pools 20

What makes it go slow? -2/2

• Not so obvious:-

- AIX tuning
- Oracle/DB2 parameters log place, SGA, Buffers
- · Read vs write characteristics
- Adapter placement, overloading bus speeds
- Throttling effects e.g., single-thread dependency
- Application errors
- Background processes (backups, batch processing) running during peak online times?
- · Concurrent access to the same files
- Changes in shared resources
- Network Buffers
- Backlevel server or I/O firmware
- · No cores left to grow beyond entitlement
- Missing AIX or VIO patches



Have a Plan What do you hope to accomplish? Describe the problem. 1. 2. Measure where you're at (baseline). Perfpmr, your own scripts,nmon, hmcscanner, etc 3. Recreate the problem while getting diagnostic data (perfpmr, your own scripts, etc.). 4. Analyze the data. 5. Document potential changes and their expected impact, then group and prioritize them. 1. Remember that one small change that only you know about can cause significant problems so document ALL changes 6. Make the changes. Group changes that go together if it makes sense to do so but don't go crazy 1. 7. Measure the results and analyze if they had the expected impact; if not, then why not? 8. Is the problem still the same? If not, return to step 1. 9. If it's the same, return to step 3. This may look like common sense but in an emergency that is the first thing to go out the window Also, find a quiet place to work so you can focus. If you are trying to work on a critical problem have someone in the team who is responsible to report back on the status so you can concentrate on the issue. Tuning is iterative – you may not get it right the first time 23



Basic Checklist for Performance Patching •Ensure your firmware is current. •Follow the memory plug-in rules. Lots of DIMMs Ensure OS level is current. Evaluate the use of SMT8. Right-size your shared LPARs. Use rPerfs for sizing VP:E ratio Java JDKs New compilers •Recompile code DO NOT USE AIX RESTRICTED TUNABLES Throughput mode Size your VIO servers correctly Tune your virtual ethernets 24

Avoiding Problems

- Stay current
- Known memory issues with 6.1 tl9 sp1 and 7.1 tl3 sp1
- Java 7.1 SR1 or higher is the preferred Java for POWER7 and POWER8
- Java 6 SR7 is minimal on POWER7 but you should go to at least Java 7
- WAS 8.5.2.2 or higher
- Refer to Section 8.3 of the Performance Optimization and Tuning Techniques Redbook SG24-8171
- HMC v9.1M940SP1 is required for POWER9 S922
 - For FW950 you need 9.1M950
 - - does not support servers prior to POWER7 (specific models and FW on P7)
 - Need to go to v9 anyway as v8 no longer supported
- Remember not all workloads run well in the shared processor pool some are better dedicated
 - Apps with polling behavior, CPU intensive apps (SAS, HPC), latency sensitive apps (think trading systems)



Resources

- Entitlement
- VPs
- Memory
- HBA queues
 - max_transfer, num_cmd_elems
- HDISK settings
 - max_transfer, queue_depth, reserve, algorithm
- Network
 - Virtual Buffers
 - TCP and UDP senda nd receive space

27

General Server Sizing thoughts · Correct amount of processor power · Balanced memory, processor and I/O · Min, desired and max settings and their effect on system overhead · Memory overhead for page tables, TCE, etc that are used by virtualization Shared or dedicated processors Capped or uncapped If uncapped – number of virtual processors Do not starve your VIO servers or client LPARs! · Set entitlement and VPs correctly · Be cautious of sizing studies - they tend to undersize memory and sometimes cores and usually do not include the VIO server needs · Consider whether the workload will play well with shared processors Never underestimate the power of common sense · Scale by rPerf (or other benchmark data) NOT by ghz when comparing boxes • Do not size the box by the sum of the entitlements. Leave room for LPARs to shrink and grow 28

CPU Sizing

When sizing the server:

Do not sum the entitlements and size the server based on entitlement The shared processor pool is supposed to be used to share processors and you need room to grow and shrink

Sum the peak concurrent VPs needed and add some for growth and then size the server

i.e.

LPAR	Entitlement	Max at noon	Max at 6pm
Α	.5	1.5	2
В	.8	1	.8
С	6	10	12
D	3	8	2
TOTAL	10.3	20.5	16.8

This system needs 21 cores not 11

29



Sample /etc/tunables/rc-tunebufs.sh

This tunes buffer settings for the two virtual adapters – assumes ent4, ent5 are virtuals

lsdev –C | grep ent will show the adapters so you can pick the right ones

```
#! /bin/ksh
#
chdev -l ent4 -a buf_mode=min -P
chdev -l ent5 -a buf_mode=min -P
chdev -l ent4 -a max_buf_tiny=4096 -P
chdev -l ent4 -a max_buf_small=4096 -P
chdev -l ent5 -a max_buf_tiny=4096 -P
chdev -l ent5 -a max_buf_small=4096 -P
chdev -l ent5 -a max_buf_small=4096 -P
```

31



HBA max_xfer_size	
The default is 0x100000 /* Default io_dma of 16MB */	
After that, 0x200000,0x400000,0x80000 gets you 128MB	
After that 0x1000000 checks for bus type, and you may get 256MB, or 128MB	
There are also some adapters that support very large max_xfer sizes which can possibly allocate	512MB
VFC adapters inherit this from the physical adapter (generally)	
Unless you are driving really large IO's, then max_xfer_size on the HBA is rarely changed beyond which provides a 128MB DMA Client setting cannot be higher than the VIOs were booted with	0x200000
	33



Entitlement and VPs Utilization calculation for CPU is different between POWER5, 6, 7, 8 and POWER9 • VPs are also unfolded sooner (at lower utilization levels than on P6 and P5) · May also see high VCSW in Iparstat · This means that starting with POWER7 you need to pay more attention to VPs You may see more cores activated a lower utilization levels But you will see higher idle · If only primary SMT threads in use then you probably have excess VPs Try to avoid this issue by: Reducing VP counts · Use realistic entitlement to VP ratios 10x or 20x is not a good idea • Try setting entitlement to .6 or .7 of VPs Per Nigel • For E=0.05 to 0.90 use VPs=1 • For E-1 to 4 round up to next whole integer • For E=4 to 8 round up 1 or 2 VPs · Ensure workloads never run consistently above 100% entitlement Too little entitlement means too many VPs will be contending for the cores • NOTE - VIO server entitlement is critical - SEAs scale by entitlement not VPs • This applies to all LPARs using the SEA as well · All VPs have to be dispatched before one can be redispatched · Performance may (in most cases, will) degrade when the number of Virtual Processors in an LPAR exceeds the number of physical processors The same applies with VPs in a shared pool LPAR - these should not exceed the cores in the pool 35

35

	Scaling	VPs afte	r Upgrade	S		
•	VPs and entitleme rPerf per core incr Example	ents should be review eases as does threa	ved when upgrading se ding depending on arc	rvers hitecture		
•	Server	Cores	smt4 rperf	rperf/core Ir	crease over p740	
	p740 p7+	16 (3.6g)	197.70	12.36		
	S822 p8	16 (4.1g)	302.40	18.90	1.53	
	S922 p9	16 (3.9max)	313.10	19.57	1.58	
•	Server	Cores	smt8 rperf	rperf/core Ir	crease over S822 SMT4	
	S822 p8	16 (4.1g)	323.60	20.23	1.07	
	S922 p9	16 (3.9max)	394.50	24.66	1.31	
TI Fo	ne above are only a or P8 and P9 SMT8 The S922 above i The S822 is 323.0	approximations using 3 gives an additional in SMT8 is 394.5 rPerf (60 (about 1.27 of the S8	published rPerf boost for threaded wor about 1.26 of the S922 SM 22 SMT4 number)	kloads IT4 number)		
rF	erf is not a guarant	tee – it is used to giv	e an idea of relative pe	rformance between	machines	
Tł lf gr	ne point to note: you are using 16 co owth or latent dema	ores (or VPs) on the and.	p740 or the S822 ther	i you do not need 1	6 cores (or VPs) on your new S922 unl	ess you have enormous
R	esize VP assignm	ents LPAR by LPAF	R, based on rperf need	ded per LPAR		
Тс	o many cores or V	Ps wastes dispatche	r and hypervisor cycles	3		
						36





Applications and SPLPARs Applications do not need to be aware of Micro-Partitioning Not all applications benefit from SPLPARs

- Applications that may not benefit from Micro-Partitioning:
 - Applications with a strong response time requirements for transactions may find Micro-Partitioning detrimental:
 - · Because virtual processors can be dispatched at various times during a timeslice
 - May result in longer response time with too many virtual processors:
 - Each virtual processor with a small entitled capacity is in effect a slower CPU
 - Compensate with more entitled capacity (2-5% PUs over plan)
 - · Applications with polling behavior
 - · CPU intensive application examples: DSS, HPC, SAS
- Applications that are good candidates for Micro-Partitioning:
 - Ones with low average CPU utilization, with high peaks:
 - Examples: OLTP, web applications, mail server, directory servers
- In general Oracle databases are fine in the shared processor pool
- · For licensing reasons you may want to use a separate pool for databases



Throughput mode in POWER9

From Steve Nasypany Session at pTechu October 2018

The VP code is aware of the core architecture and will place/collapse smaller workloads slightly more aggressively when workloads are present

• Optimization has the additional impact of reducing physical consumption

• Single thread utilization is calibrated to ~32% in SMT8 (~44% in SMT4), so below the default VP dispatch threshold of ~50% per core

• Because single-threads are calibrated lower and equivalent workloads will overall generate a lower utilization, they are more likely to fall below dispatch threshold, thus lowering physical consumed

• Linux not running on PowerVM dispatches to SMT4 cores and does not use hardware calibration











%user	%sys	%wait	%idle		
16.8	 28.7	6.4	48.1		
17.0	29.3	5.8	48.0		
busy = % parstat - System o	and smt %occupati h 30 2 ou configurat	=4 means ion of the L itput ion: type=[I have 80/ _CPUs at t	/ 4=20 real he system mode=Cap	n and user level apped smt=4 lcpu=80 mem=524288MB
lbusy = % parstat - System o %user	and smt %occupat h 30 2 ou configurat %sys	=4 means ion of the L itput ion: type=I %wait	L have 80/ CPUs at t Dedicated %idle	/ 4=20 real he system mode=Cap %hypv	n and user level apped smt=4 lcpu=80 mem=524288MB hcalls
logu=80 parstat - System o %user 16.8	and smt %occupat h 30 2 ou configurat %sys 29.8	=4 means ion of the L itput iion: type=I %wait 5.4	I have 80/ CPUs at t Dedicated %idle 48.0	/4=20 real he system mode=Cap %hypv 61.3	apped smt=4 lcpu=80 mem=524288MB hcalls 2222545



b814aix1: lj	parstat	30 2																	
System con	figuratio	on: type	e=Share	ed mode	=Uncap	oed sm	t=8 lc	ou=48	mem=	32768	MB psiz	e=2 ent	=0.50						
%user		%sys	%wait	%idle	physc	%entc	lbusy	app	vcsw	phint									
0.0 0.0		0.1 0.2	0.0 0.0	99.9 99.8	0.00 0.00	0.8 1.0	2.3 2.3	1.96 1.96	244 257	0 0									
b814aix1: n	npstat -	·s																	
System con	figuratio	on: Icpu	u=48 en	t=0.5 m	ode=Uno	apped													
	P (²roc0 0.00%									F	Proc8 .00%							
cpu0 cp 0.00% 0.0	u1 c)0% 0.	;pu2 .00%	cpu3 0.00%	cpu4 0.00%	cpu5 0.00%	cpu6 0.00%	cr (ou7).00%	C 0	pu8 .00%	cpu9 0.00%	cpu10 0.00%	cpu11 0.00%	cpu12 0.00%	cpu13 0.00%	cpu14 0.00%	cpu15 0.00%		
	Pi (roc16 0.00%									F	roc24 .00%							
cpu16 cpu 0.00% 0.0	17 cpu)0% 0.	u18 cj .00% (pu19 0 0.00%	cpu20 (0.00% (0.00% (pu22).00%	cpu23 0.009	8 %	0 0	ou24 .00%	cpu25 0.00%	cpu26 0.00%	cpu27 0.00%	cpu28 0.00%	cpu29 0.00%	cpu30 0.00%	cpu31 0.00%		
	Pi (roc32 0.00%									F	roc40 .00%							
cpu32 cp	J33 cp	u34	cpu35	cpu36	cpu37	cpu38	cpu 6 0 0	139 0%	c c	pu40	cpu41 0.00%	cpu42	cpu43	cpu44	cpu45	cpu46	cpu47		

vmstat -IW bnim: vmstat -IW 2 2 vmstat -IW 60 2 System configuration: lcpu=12 mem=24832MB ent=2.00 memory kthr page faults CDU _____ ----- ----fre fi fo pi po fr sr in rbpw avm r b p w avm fre fi fo pi po fr sr in sy cs us sy id wa pc ec 3 1 0 2 2708633 2554878 0 46 0 0 0 0 3920 143515 10131 26 44 30 0 2.24 112.2 sy cs us sy id wa pc 6 1 0 4 2831669 2414985 348 28 0 0 0 0 2983 188837 8316 <u>38</u> <u>39</u> 22 0 **2.42 120.9** Note pc=2.42 is 120.0% of entitlement When looking at system time to user time ratios - remember on a VIO server that high system time is most likely normal as the VIO handles all the I/O and network and really has little normal user type work -I shows I/O oriented view and adds in the p column p column is number of threads waiting for I/O messages to raw devices. -W adds the w column (only valid with -I as well) w column is the number of threads waiting for filesystem direct I/O (DIO) and concurrent I/O (CIO) r column is average number of runnable threads (ready but waiting to run + those running) This is the global run queue - use mpstat and look at the rq field to get the run queue for each logical CPU b column is average number of threads placed in the VMM wait queue (awaiting resources or I/O) 4٩



topas -C

Topas CEC 1	Monit	tor				Int	erva	al:	1	0		Thu Fe	b 27	08:5	3:05 2014
Partitions	Memo	ory	(GB)				Pro	ces	sor	3	_				
Shr: 5	Mon	:86.	0 II	nUse	:23.	.0	Shr	: 1	В	PSz: 16	Don	: 0.0 S	hr Ph	узВ	0.02
Ded: 0	Avl	:	1			1	Ded	: (0	APP: 16.0	0 Stl	: 0.0 D	ed_Ph	ysB	0.00
Host	OS	Mod	Mem	InU	Lp	Us	Sy	Wa	Id	PhysB	Vcsw	Ent 🖁	EntC	PhI	pmem
								-sha	are	d					
b740n11	A71	Ued	32	5.3	32	0	0	0	99	0.03	210	4.00	0.7	0	-
b740vio2	A61	U-d	3.0	2.8	8	0	0	0	99	0.00	256	0.50	0.8	0	-
b740ft2	A71	Ued	32	5.3	4	0	0	0	99	0.00	191	1.00	0.4	0	-
	A61	U-d	3.0	2.8	4	0	0	0	99	0.00	171	0.50	0.6	0	-
b74011	A71	U-d	16	7.1	16	0	0	0	99	0.00	212	2.00	0.1	0	
Host	OS	Mod	Mem	InU	Lp	Us	Sy	Wa	Id	PhysB	Vcsw	&istl	<pre>%bst</pre>	1	

Shows pool size of 16 with all 16 available Monitor VCSW as potential sign of insufficient entitlement

51

Lotopas_nmonqqC=many-CPUsqqqqqqqdbst=b740ft2qqqqqqqqRefresh=2 se	csqqq09:09.31qqqqqqqqqqqqqqqqqqqqqqqqqqqqq
» Shared-CPU-Logical-Partition qqqqqqqqqqqqqqqqqqqqqqqqqqqqqqqqqqqq	aadaaadaaadaaadaaadaaadaaadaaaaa
xPartition:Number=14 "b740ft2"	
xriags:LFARed DRable SMI Shared Uncapped FoolAuth Migratable Not-	Donating AMSable.
xSummary: Entitled= 1.00 03ed 0.00 (0.5%) 0.0% of CPUs in	Pool
x POLICIUS=16 Unused 15.96 U.Us OI CPUS II.	. F001
wmay Dhug in gug 16 Can Drocegor Min 0 10 SDIDAD Group. Dool	32782+0
when rules in Sys 16 Cap. Processor Max 4 00 Memory (MR) Min Ma	v 1024:65536
Wintual Online 1 Cap. Increment. 0.01 Memory(MB) Online	32768
xLogical Online 4 Cap. Unallocated 0.00 Memory Region LME	256MB min
xPhysical pool 16 Cap. Entitled 1.00 Time	Seconds
xSMT threads/CPU 4 -MinRegVirtualCPU 0.10 Time Dispatch Whe	el 0.0100
xCPU Min-Max Weight MaxDispatch Later	cy 0.0000
xVirtual 1 4 Weight Variable 128 Time Pool Idle	15.9761
xLogical 1 16 Weight Unallocated 0 Time Total Dispat	ch 0.0046
x xEvent= 0 SerialNo Old= Current=F6934C When=	
xx xShared_Pools MaxPoolCapacity=16.00 MyPoolMax =16.00 SharedCPU-	 Total=16.00
xSharedCPU=16 EntPoolCapacity=15.20 MyPoolBusy= 0.02 SharedCPU-	Busy = 0.02
xaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaa	adaddaddadadadadadadadadadadadada
X	
v	



emstat	
emstat -a 30 2	
Alignment Alignment Emulation Emulation SinceBoot Delta SinceBoot Delta 0 0 0 0 0 0 0 0 Check for emulations and alignments Caused when old code run on new architectures or operating systems Code should be recompiled whenever possible when upgrades occur	
Binary compatibility statement means it will work but the best performance is obtained wi version compiled for the architecture	hen you run a
	54

56



7b. MEMORY

55

Memory Types

- Kernel
 - This is AIX itself and the minimum is around 2GB.
- Persistent
 - · Backed by filesystems (called file cache)
 - Raw devices, GPFS, Oracle ASM does not make use of file cache.
- Working storage
 - Dynamic
 - · Includes executables and their work areas
 - · Backed by page space
 - Shows as avm in a vmstat –I (multiply by 4096 to get bytes instead of pages) or as %comp in nmon analyser or as a percentage of memory used for computational pages in vmstat –v
 ALSO NOTE if %comp is near or >97% then you will be paging and need more memory
- Prefer to steal from persistent as it is cheap
- minperm, maxperm, maxclient, lru_file_repage and page_steal_method all impact these decisions















Affinity

- LOCAL SRAD, within the same chip, shows as s3
- NEAR SRAD, within the same node intra-node, shows as s4
- FAR SRAD, on another node inter-node, shows as s5
- Command is Issrad -av or can look at mpstat -d
- Topas M option shows them as Localdisp%, Neardisp%, Fardisp%
- The further the distance the longer the latency
- · Problems you may see
 - SRAD has CPUs but no memory or vice-versa
 - CPU or memory unbalanced
- Note on single node systems far dispatches are not as concerning
- To correct look at new firmware, entitlements and LPAR memory sizing
- Can also look at Dynamic Platform Optimizer (DPO)

64

Memory Tips

Avoid having chips without DIMMs. Attempt to fill every chip's DIMM slots, activating as needed. Hypervisor tends to avoid activating cores without "local" memory.



Diagram courtesy of IBM

65

mpsta	at -dw 1	5 2 ou	tput													
Syste	em confi	gurati	on: lcpu	=80 mo	de=Cappe	ed							lc	ocal	near	far
cpu	cs	ics	bound	\mathbf{rq}	push S	3pull	S3grd	SOrd	Slrd	S2rd	S3rd	S4rd	S5rd S	33hrd	S4hrd	S5hrd
0	92	31	1	1	0	0	44	95.1	0.6	0.0	1.2	0.0	3.1 1	100.0	0.0	0.0
1	3	0	0	0	0	0	0	82.1	12.0	0.0	6.0	0.0	0.0 1	100.0	0.0	0.0
2	2	0	0	0	0	0	0	95.3	4.7	0.0	0.0	0.0	0.0 1	100.0	0.0	0.0
3	2	0	0	0	0	0	0	95.1	4.9	0.0	0.0	0.0	0.0 1	100.0	0.0	0.0
4	314	21	1	1	0	0	95	97.4	0.1	0.0	0.6	0.0	2.0 1	100.0	0.0	0.0
5	14	2	0	0	0	0	0	94.5	5.5	0.0	0.0	0.0	0.0 1	100.0	0.0	0.0
6	3	0	0	0	0	0	0	96.4	3.6	0.0	0.0	0.0	0.0 1	100.0	0.0	0.0
7	2	0	0	0	0	0	0	95.3	4.7	0.0	0.0	0.0	0.0 1	100.0	0.0	0.0
			lots :	more 1	ines her	e										
72	280	38	1	1	0	224	32	91.0	0.0	0.0	3.0	0.0	6.0	66.5	8.6	24.8
73	2	0	0	0	0	0	0	95.6	4.4	0.0	0.0	0.0	0.0	95.6	0.0	4.4
74	2	0	0	0	0	0	0	95.1	4.9	0.0	0.0	0.0	0.0	95.1	0.0	4.9
75	2	0	0	0	0	0	0	94.9	5.1	0.0	0.0	0.0	0.0	94.9	0.0	5.1
76	161	25	0	0	0	87	14	93.0	0.0	0.0	2.9	0.0	4.1	77.1	7.4	15.4
77	2	0	0	0	0	0	0	94.3	5.7	0.0	0.0	0.0	0.0	94.3	0.0	5.7
78	1	0	0	0	0	0	0	93.8	6.2	0.0	0.0	0.0	0.0	93.8	0.0	6.2
79	2	0	0	0	0	0	0	92.3	7.7	0.0	0.0	0.0	0.0	92.3	2.6	5.1
ALL	6771	996	13	13	0	1888	2170	93.8	0.3	0.0	2.0	0.0	3.9	89.5	3.5	6.9

			•														
Syst	em conf	igurati	ion: lcpu	u=120 m	ode=Ca	oped							loca	al nea	ır far		
		ice	hound		nuch	62mu11	62 and	cond	C1nd	cand	cand	C And	SEnd Sohn	d Sábrd	l CEhnd	%non	
cpu a	6727	220	Dound	19	pusn	2100	102/	20 E	511.0	521°U	351°u	2 0	551°U 55111°	u 54m.u		/₀nsp 101	
1	5766	265	4 0	6	9	1874	1634	67 7	10 7	0.0	17.6	1.9	0.0 83.	8 18 7	0.0	101	
2	5441	265	0 0	â	â	10/4	1034	73.1	10.6	0.0	16.4	0.0	0.0 95	2 4.8	0.0	101	
3	5667	292	ø	õ	0	ø	õ	73.6	10.1	0.0	16.3	0.0	0.0 95.	1 4.9	0.0	101	
4	0	0	0	0	0	0	ø	0.0	100.0	0.0	0.0	0.0	0.0 88.	9 11.1	0.0	101	
5	0	0	0	0	0	0	0	0.0	100.0	0.0	0.0	0.0	0.0 100.	ø ø.e	0.0	101	
6	0	0	0	0	0	0	0	0.0	100.0	0.0	0.0	0.0	0.0 100.	ø ø.e	0.0	101	
7	0	0	0	0	0	0	0	0.0	100.0	0.0	0.0	0.0	0.0 71.	4 28.6	0.0	101	
8	5895	298	1	11	0	1974	1654	67.5	11.0	0.0	17.5	4.1	0.0 81.	5 18.5	0.0	101	
9	5758	291	0	0	0	1948	1600	67.7	10.7	0.0	17.5	4.1	0.0 81.	9 18.1	0.0	101	
10	5242	235	0	0	0	0	0	71.2	11.9	0.0	17.0	0.0	0.0 94.	6 5.4	0.0	101	
11	6319	332	1	9	0	0	0	74.1	10.0	0.0	15.9	0.0	0.0 94.	9 5.1	0.0	101	
12	0	0	0	0	0	0	0	0.0	100.0	0.0	0.0	0.0	0.0 100.	0 0.0	0.0	101	
13	0	0	0	0	0	0	0	0.0	100.0	0.0	0.0	0.0	0.0 71.	4 28.6	6.0	101	
14	0	0	0	0	0	0	0	0.0	100.0	0.0	0.0	0.0	0.0 100.	0 0.0	0.0	101	
15	0	0	0	0	0	0	0	46.2	53.8	0.0	0.0	0.0	0.0 76.	9 23.1	0.0	101	
118	0	0	0	0	0	0	0	0.0	100.0	0.0	0.0	0.0	0.0 66.	7 33.3	3 0.0	101	
119	0	0	0	0	0	0	0	0.0	100.0	0.0	0.0	0.0	0.0 50.	0 50.0	0.0	101	
	265322	12290	43	112	0	60475	99151	72.1	10.0	0.0	13.9	4.0	0.0 86.	7 13.3	8 0.0	0	



vmstat –v Output

memory pools	4	
numperm	10.1%	
numclient	10.1%	
uptime	up 16 davs.	4:46

17308 pending disk I/Os blocked with no pbuf 4538737 paging space I/Os blocked with no psbuf 1972 file system I/Os blocked with no fsbuf 8724 client file system I/Os blocked with no fsbuf 7077 external pager file system I/Os blocked with no fsbuf 89.9 percentage of memory used for computational pages

pbufs pagespace JFS NFS/VxFS JFS2

69

numclient=numperm so most likely the I/O being done is JFS2 or NFS or VxFS Based on the blocked I/Os it is clearly a system using JFS2 It is also having paging problems (vmstat and lsps commands) pbufs also need reviewing (lvmo command)

Also note that computational was 89.9%

69



Calculating minfree and maxfree vmstat –v | grep memory 3 memory pools vmo -a | grep free maxfree = 1088 minfree = 960Calculation is: minfree = (max (960,(120 * lcpus) / memory pools)) maxfree = minfree + (Max(maxpgahead,j2_maxPageReadahead) * lcpus) / memory pools So if I have the following: Memory pools = 3 (from vmo - a or kdb) J2_maxPageReadahead = 128 CPUS = 6 and SMT on so lcpu = 12 So minfree = (max(960,(120 * 12)/3)) = 1440 / 3 = 480 or 960 whichever is larger And maxfree = minfree + (128 * 12) / 3 = 960 + 512 = 1472 I would probably bump this to 1536 rather than using 1472 (nice power of 2) The difference between minfree and maxfree should be no more than 1K per IBM If you over allocate these values it is possible that you will see high values in the "fre" column of a vmstat and yet you will be paging.





POWER9 Performance Best Practices

A brief checklist

This document is intended as a short summary for customers on key items that should be looked at when planning a migration. For a more in-depth and more complete set of recommendations, please refer to the document links provided on the second page
 Instructions

 Fix Central provides latest updates. Latest F/W levels as of this writing : F/W404 for ALL POWERS systems.

 Use the FLRT too to obtain the recommended levels for a given platform.

 NOTE: Ensure required HAC level is installed whom updating F/W.

 AII POWER: Follow proper memory plug-in rules Fix Central provides the latest updates for AIX, [BMI, VIOS, Linux, HMC and F/W. In addition to that, the FLRT tool provides the recommended levels for each H/W model. Use these tools to mainfain your system us to date age of the performance of FOWERR CPUs, we recommed leither evaluate the use of SMTB. Proper sizing is also recommended to maximize the FOWERR improvement. We recommend when moving to SMTB to reboot the partition to gath the best performance of this charge.

 RHE[27: For network bandwidth sensitive workloads, we recommend increase the receive queue size from 1024 to 8122.
 Instructions
 VIOS 3.1 runs on native POWER9 mode Description Description Ensure firmware is VIOS 3.1 runs on native POWER9 mode
 VIOS 2.2.5 and above provide performance improvements over prior versions.
 Use of 10Gb dedicated network for LPM is preferred. urrent - use of user version terrorium ret LPM Is pretered.
 Use of dual MSPs can improve performance (minimum requirements, VIOS 2.2.5 and F/W860)
 Tuning a VIOS is not recommended unless directed by VIOS/AIX support.
 Restricted transless should not be modified (unless directed by AIX/VIOS development)
 Tunables should not be migrated across H/W or AIX levels. Memory DIMMs Ensure OS level current AIX Tunables / VIOS Tunables Tunables should not be migrated across HW or AIX levels.
 The AIX OS system is optimized for best raw throughput at higher CPU usage. If the customer requires to reduce CPU usage (cr), use the schedu tanable ymm, throughput_mode to tune the workload and evaluate the benefits of raw throughput s. CPU usage.
 Shared Ethermer, daagters using a 1053, 4050 er 0006 to HU ether workload and evaluate the benefits of raw throughput is, choice and the schedu tanable of the two workload and evaluate the benefits of raw throughput is, a CPU usage.
 Shared Ethermer, daagters using a 1053, 4050 er 0006 bit Helmox adapter as a backing device should enable the "file_n_run" attribute, wis choice, on the network adapter port i.
 If configured with hained processors:
 Assign total entitlement of all VIOS partitions to be 10-15% of cores in shared pool and assign CPU ratio of 2.1 (CPU table). Refer to the PowerVIM Best Practices for additional recommendations
 Assign uncapped mode and set variable weight capacity of VIOS partition higher than all client LPARs serviced by VIOS
 For vFC, ensure no more than 64 client connections tabl per physical fits adapter on the VIOS. Also, ensure no more than 64 client connections table per physical fits adapter on the VIOS.
 For vFC, ensure no more than 64 client connections table per typical fits adapter on the VIOS. Also, ensure no more than 64 client connections table per physical fits adapter on the VIOS.
 For vFC, ensure no more than 64 client connections table per typical fits adapter on the VIOS. Also, ensure no more than 64 client connections table sets than or equal the queue_depth of the physical disk in the VIOS.
 For vFC, ensure no more than 64 client connections table sets than or equal the queue_depth of the physical disk in the VIOS.
 For vFC, ensure no more than 64 client connections table to the ticlient. These are physical limits, pa SMT8 AIX CPU utilization 40GbE adapter VIOS configuratio om 1024 to 8192 from 1024 to 8192.
When migrating to POWER9, we recommend considering using SMT8, and size the LPARs based the SMT8 Perf values; in many instances, this will likely reduce the number of VPs required. Use Workload Estimator (WLE) for sizing LPARs for CPU consumption as it provides better sizing ing a syster Use Workload Estimator (WLE) for sizing LPARs for CPU consumption as It provides better sizing results.
Assign entitled capacity (EC) to sustained peak utilization for LPARs with critical SLA requirements.
Assign Etc varenge utilization and number of virtual CPUs to peak utilization(s)visical core consumption) for LPARs with non-critical SLA.
Ensure the average LLRA utilization is equal or less than 75% of the entitled capacity.
Current FW evide ensure optimal placement of the partitions. However, if constant DLPAR operations current level of FM).
Use the CPL is the commended the use DPO to optimite placement (requires current level of FM).
Is DLA USE TO A CPL in the CPL is a recommended the use DPO to optimite placement (requires current level of AL).
Is DLA USE TO A CPL in the CPL is the commended the use DPO to optimite placement (requires the current level of AL).
Is DLA USE TO A CPL in the CPL is the commended the use DPO to optimite placement (requires the current level of AL).
Is DLA USE TO A CPL in the CPL is the commended the use DPO to optimite placement (requires the current level of AL).
Is DLA USE TO A CPL in the CPL is the commended the use DPO to optimite placement (requires the current level of AL).
Is DLA USE TO A CPL is the CPL is the commended the use DPO to optimite placement (requires the current level of AL).
Is DLA USE TO A CPL is DLA USE TO A CPL is the CPL is Right-size your Shared LPARs Partition Placement IBM XL CC+ For AIX V16.1 and XL Fortran V16.1 added support for POWER9. Also adds support for C++11 and C++14.
 IBM XL CC+ For Linux V16.1.1 & XL Fortran V16.1.1 support for P9 ISA Advanced Toolchain: 11.0.3 and later gcc: Version 7 gc ce is recommended for P9 ISA support. Also includes support for "-mtune=power9" AIX: Change the following vNIC interface settings as follows: chefev - lente⁻ arc. rque_rum® - a tr_que_rum® - arc_que_elem=2048 - a tx_que_elem=1024 - a use_rec_q_ruleno Linux distors - Update to listest kernel. Also, set rr/tx queues to maximum method - Lenter & a tr_que_ruleno IBM JCK 678 & the minimum level to exploit POWER9 - Open JDK 12 provides partial support for P9 ISA - User gC4R ste pages normally increases application performance Compiler 2) (virtual_q_depth - 3) Only enable the largeend attribute on the SEA (physical adapter backing the SEA) if all LPARs serviced by the VIOS are AIX partitions. Increase the virtual Etherent (eth) device driver buffers if the partition is dropping packets on the virtual interface even when running with entitled CPU capacity. e.g., chdev = entit - an ax_budg vocenNINN NOTE: For desired buffer size adjustments, refer to 'AIX on Power – Performance FAQ* link below Set largeend on VETH adapter to improve performance (AIX); VNIC Virtual Ethernet adapters on AIX Java Open JDK 1.8 provides partial support for P9 ISA Use of 64k size pages normally increases application performance isure Technology Updates are current (see link below) chdev -l en# -a mtu_bypass=on (or) ifconfig en# largesend Ē IBMi

Above at: http://www14.software.ibm.com/webapp/set2/sas/f/best/power9_performance_best_practices.pdf

73



Setup NTP

#vi /etc/ntp.conf Comment out broadcast and add: server 0.pool.ntp.org server 1.pool.ntp.org

#ntpdate 0.pool.ntp.org

Update rc.tcpip to start ntp at boot Now start NTP #startsrc -a "-c /home/padmin/config/ntp.conf" -s xntpd

You can substitute your own NTP servers for the ones above if you have them



Update I/O Firmware May also need to update HMC, server firmware, VIO first – use FLRTVC to check and read the readm	es
As root run: Ismcode –A	
Check on Fix Central under Power I/O Firmware You will need to know what kind of adapters you have (feature codes)	
If you are updating the primary you can let it failover or you can force a failover chdev -I ent14 -a ha_mode=standby When done chdev -I ent14 -a ha_mode=auto	
Example updating a 5899 network adapter with code uploaded to server	
cd /software/adapters/5899 rpm -ivhignoreos e414571614102004.10240310.aix.rpm diag -T download -d ent0 Updated all 4 ent0-ent3 You may have to unconfigure the SEA to do this (see next slide)	77



10. Don't use restricted tunables unless IBM support tells you to





Useful Links

- Jaqui Lynch Articles
 - <u>http://www.circle4.com/jaqui/eserver.html</u>
- Jaqui's Movie Replays
 http://www.circle4.com/movies
- Jaqui's Youtube Channel
 <u>https://www.youtube.com/user/adespota4/</u>
- Nigel Griffiths AIXpert Blog
 <u>https://www.ibm.com/support/pages/aixpert-blog-nigel-griffiths-mrnmon</u>
- Nigel Griffiths Twitter mr_nmon
 https://twitter.com/mr_nmon
- Nigel Griffiths YouTube
 - https://www.youtube.com/nigelargriffiths
- Gareth Coates Tricks of the POWER Masters

 <u>https://www.ibm.com/support/pages/node/1116939</u>

 Gareth Coates Twitter power_gaz
 - <u>https://twitter.com/power_gaz</u>





Backup Slides



Terminology Run QueueRun queue length is another well known metric of CPU usage It refers to the number of software threads that are ready to run, but have to wait because the CPU(s) is/are busy or waiting on interrupts The length is sometimes used as a measure of health, and long run queues usually mean worse performance, but many workloads can vary dramatically. It is quite possible for a pair of single-threaded workloads to contend for a single physical resource (batch, low run queue, bad performance) while dozens of multi-threaded workloads share it (OLTP, high run queue, good performance)

Terminology

Context Switches

- Context switches
 - The number of times a running entity was stopped and replaced by another
 - Collected for Threads (operating system) and Virtual Processors (hypervisor)
 - There are voluntary and involuntary context switches
- How Many "context switches" are Too Many?
 - No rules of thumb exist
 - Voluntary: Not an issue because it means no work for the CPU
 - Involuntary: Could be an issue, but generally the bottleneck will materialize in a easier to diagnosis metric; such as, CPU utilization, physical consumption, entitlement consumed, run queue
 - Establish a baseline and compare when system encounter performance problems
- Tool outputs
 - vmstat reports total context switches as "cs"
 - sar -w as "cswch/s"
 - AIX lparstat reports virtual processor context switches as "vcsw"
 - AIX mpstat reports per logical processor context switches (involuntary: *ilcs* and voluntary: *vlcs*). The sum of vlcs should roughly equal VP vcsw.

85

