



# **NETWORK TUNING in AIX**

See article at: <u>http://www.ibmsystemsmag.com/aix/administrator/networks/network\_tuning/</u> Replay at: <u>http://www.circle4.com/movies/</u>

# Tunables

- The tcp\_recvspace tunable
  - The *tcp\_recvspace* tunable specifies how many bytes of data the receiving system can buffer in the kernel on the receiving sockets queue.
- The tcp\_sendspace tunable
  - The *tcp\_sendspace* tunable specifies how much data the sending application can buffer in the kernel before the application is blocked on a send call.

## The rfc1323 tunable

3

- The *rfc1323* tunable enables the TCP window scaling option.
- By default TCP has a 16 bit limit to use for window size which limits it to 65536 bytes. Setting this to 1 allows for much larger sizes (max is 4GB)

## • The sb\_max tunable

• The *sb\_max* tunable sets an upper limit on the number of socket buffers queued to an individual socket, which controls how much buffer space is consumed by buffers that are queued to a sender's socket or to a receiver's socket. *The tcp\_sendspace attribute must specify a socket buffer size less than or equal to the setting of the sb\_max attribute* 

# **UDP Send and Receive**

## udp\_sendspace

Set this parameter to 65536, which is large enough to handle the largest possible UDP packet. There is no advantage to setting this value larger

## udp\_recvspace

5

Controls the amount of space for incoming data that is queued on each UDP socket. Once the *udp\_recvspace* limit is reached for a socket, incoming packets are discarded.

Set this value high as multiple UDP datagrams could arrive and have to wait on a socket for the application to read them. If too low packets are discarded and sender has to retransmit.

Suggested starting value for *udp\_recvspace* is 10 times the value of *udp\_sendspace*, because UDP may not be able to pass a packet to the application before another one arrives.

# Some definitions

- TCP large send offload
  - Allows AIX TCP to build a TCP message up to 64KB long and send It in one call down the stack. The adapter resegments into multiple packets that are sent as either 1500 byte or 9000 byte (jumbo) frames.
  - Without this it takes 44 calls (if MTU 1500) to send 64KB data. With this set it takes 1 call. Reduces CPU. Can reduce network CPU up to 60-75%.
  - It is enabled by default on 10Gb adapters but not on VE or SEA.
- TCP large receive offload
  - Works by aggregating incoming packets from a single stream into a larger buffer before passing up the network stack. Can improve network performance and reduce CPU overhead.
- TCP Checksum Offload
  - Enables the adapter to compute the checksum for transmit and receive. Offloads CPU by between 5 and 15% depending on MTU size and adapter.

# Large Send and Large Receive

- Important note
  - Do not enable these on the sea if used by Linux or IBM I client partitions
  - Do not enable if used by AIX partitions set up for IP forwarding

# Some more definitions

## MTU Size

- The use of large MTU sizes allows the operating system to send fewer packets of a larger size to reach the same network throughput. The larger packets greatly reduce the processing required in the operating system, assuming the workload allows large messages to be sent. If the workload is only sending small messages, then the larger MTU size will not help. Choice is 1500 or 9000 (jumbo frames). Do not change this without talking to your network team.
- MSS Maximum Segment Size
  - The largest amount of data, specified in bytes, that a computer or communications device can handle in a single, unfragmented piece.
  - The number of bytes in the data segment and the header must add up to less than the number of bytes in the maximum transmission unit (MTU).
- Computers negotiate MTU size
  - Typical MTU size in TCP for a home computer Internet connection is either 576 or 1500 bytes. Headers are 40 bytes long; the MSS is equal to the difference, either 536 or 1460 bytes.

# More on MTU and MSS

Routed data must pass through multiple gateway routers.

We want each data segment to pass through every router without being fragmented.

If the data segment size is too large for any of the routers through which the data passes, the oversize segment(s) are fragmented.

This slows down the connection speed and the slowdown can be dramatic.

Fragmentation can be minimized by keeping the MSS as small as reasonably possible.

# Adapter Options and Defaults

Table 7. Adapters and their available options, and system default settings

Adapter type	Feature code	TCP checksum offload	Default setting	TCP large send	Default setting
GigE, PCI, SX & TX	2969, 2975	Yes	OFF	Yes	OFF
GigE, PCI-X, SX and TX	5700, 5701	Yes	ON	Yes	ON
GigE dual port PCI-X, TX and SX	5706, 5707	Yes	ON	Yes	ON
10 GigE PCI-X LR and SR	5718, 5719	Yes	ON	Yes	ON
10/100 Ethernet	4962	Yes	ON	Yes	OFF
ATM 155, UTP & MMF	4953, 4957	Yes (transmit only)	ON	No	N/A
ATM 622, MMF	2946	Yes	ON	No	N/A

10

# Starter set of tunables 3

Typically we set the following for both versions:

NETWORK no -p -o rfc1323=1 no -p -o tcp\_sendspace=262144 no -p -o tcp\_recvspace=262144 no -p -o udp\_sendspace=65536 no -p -o udp\_recvspace=655360

Also check the actual NIC interfaces and make sure they are set to at least these values You can't set udp\_sendspace > 65536 as IP has an upper limit of 65536 bytes per packet

Check sb\_max is at least 1040000 - increase as needed



	My VIO Server SEA
# ifconfi	g -a
en6: flags=1e	080863 580<119 BROADCAST NOTRAILERS BLINNING SIMPLEX MULTICAST GROUPRT 6
4BIT,CHI	ECKSUM_OFFLOAD(ACTIVE),CHAIN>
inet	192.168.2.5 netmask 0xffffff00 broadcast 192.168.2.255
tcp	_sendspace 262144 tcp_recvspace 262144 rfc1323 1
lo0:	
flags=e0	8084b,1c0 <up,broadcast,loopback,running,simplex,multicast,grouprt,64bi< td=""></up,broadcast,loopback,running,simplex,multicast,grouprt,64bi<>
inet	: 127.0.0.1 netmask 0xff000000 broadcast 127.255.255.255
inet	6 ::1%1/0
icp	_sendspace 1510/2 (cp_1ecvspace 1510/2 (iC1525 1
13	

Interface	Speed	MTU	tcp_sendspace	tcp_recvspace	rfc1323	tcp_nodelay	tcp_mssdflt
o0 (loopback)	N/A	16896	131072	131072	1		
Ethernet	10 or 100 (Mbit)						
Ethernet	1000 (Gigabit)	1500	131072	65536	1		
Ethernet	1000 (Gigabit)	9000	262144	131072	1		
Ethernet	10 GigE	1500	262144	262144	1		
Ethernet	10 GigE	9000	262144	262144	1		
Ether Channel	Configures based	on speed	d/MTU of the und	lerlying interfaces	š.	•	
Virtual Ethernet	N/A	any	262144	262144	1		
nfiniBand	N/A	2044	131072	131072	1		
	en from AIX V7.1 Pe	rrormanc	e luning Guide				
Check up t Aix V5.3	o date information w-01.ibm.com/sup	at: port/knov	vledgecenter/api/o	content/ssw_aix_5	i3/com.ibm	.aix.prftungd/doo	:/prftungd/prftungd_pdf.p





# Network Performance and Throughput

## • Depends on:

- Available CPU power entitlement at send/receive VIOs and client LPARs
- MTU size
- Distance between receiver and sender
- Offloading features
- Coalescing and aggregation features
- TCP configuration
- Firmware on adapters and server
- Ensuring all known efixes are on for 10GbE issues
- Pay attention to adapter type and placement
- Use Isslot –c pci
  - This helps you figure out what kind of slots you have

Network type	Raw bit rate (Mbits)	Payload rate (Mb)	Payload rate (MB)
10 Mb Ethernet, Half Duplex	10	5.8	0.7
10 Mb Ethernet, Full Duplex	10 (20 Mb full duplex)	18	2.2
100 Mb Ethernet, Half Duplex	100	58	7.0
100 Mb Ethernet, Full Duplex	100 (200 Mb full duplex)	177	21.1
1000 Mb Ethernet, Full Duplex, MTU 1500	1000 (2000 Mb full duplex)	1811 (1667 peak) 1	215 (222 peak) <sup>1</sup>
1000 Mb Ethernet, Full Duplex, MTU 9000	1000 (2000 Mb full duplex)	1936 (1938 peak) 1	231 (231 peak) 1
10 Gb Ethernet, Full Duplex, MTU 1500	10000 (20000 Mb full duplex)	14400 (18448 peak) 1	1716 (2200 peak) 1
10 Gb Ethernet, Full Duplex, MTU 9000	10000 (20000 Mb full duplex)	18000 (19555 peak) 1	2162 (2331 peak) 1
FDDI, MTU 4352 (default)	100	97	11.6
ATM 155, MTU 1500	155 (310 Mb full duplex)	180	21.5
ATM 155, MTU 9180 (default)	155 (310 Mb full duplex)	236	28.2
ATM 622, MTU 1500	622 (1244 Mb full duplex)	476	56.7
ATM 622, MTU 9180 (default)	622 (1244 Mb full duplex)	884	105

# Notes on 10GbE

- Using jumbo frames better allows you to use the full bandwidth coordinate with network team first
  - Jumbo frames means an MTU size of 9000
  - Reduces CPU time needed to forward packets larger than 1500 bytes
  - Has no impact on packets smaller than 1500 bytes
  - Must be implemented end to end including virtual Ethernet, SEAs, etherchannels, physical adapters, switches,
  - core switches and routers and even firewalls or you will find they fragment your packets
  - Throughput can improve by as much as 3X on a virtual ethernet
- Manage expectations
  - Going from 1GbE to 10GbE does not mean 10x performance
  - You will need new cables
  - You will use more CPU and memory
    - Network traffic gets buffered
    - This applies to the SEA in the VIOS
- Check that the switch can handle all the ports running at 10Gb
- Make sure the server actually has enough gas to deliver the data to the network at 10Gb

19

# 10GbE Tips

- Use flow control everywhere this stops the need for retransmissions
  - Need to turn it on at the network switch,
  - Turn it on for the adapter in the server
    - chdev –l ent? –a flow\_cntrl=yes
- If you need significant bandwidth then dedicate the adapter to the LPAR
   There are ways to still make LPM work using scripts to temporarily remove the adapter
- TCP Offload settings largesend and large\_receive
  - These improve throughput through the TCP stack
- Set largesend on (TCP segmentation offload) should be enabled by default on a 10GbE SR adapter
  - AIX chdev –l en? –a largesend=on
  - On vio chdev –dev ent? –attr largesend=1
  - With AIX v7 tl1 or v6 tl7 chdev –l en? –l mtu\_bypass=on
- Mtu\_bypass
  - At 6.1 tl7 sp1 and 7.1 sp1
  - O/s now supports mtu\_bypass as an option for the SEA to provide a persistent way to enable largesend
  - See section 9.11 of the AIX on POWER Performance Guide

# 10GbE Tips

- Try setting large\_receive on as well (TCP segment aggregation)
  - AIX chdev –I en? –a large\_receive=on
  - VIO chdev –dev ent? –attr large\_receive=1
- If you set large\_receive on the SEA the AIX LPARs will inherit the setting
- Consider increasing the MTU size (talk to the network team first) this increases the size of the actual packets
  - chdev -l en? mtu=65535 (9000 is what we refer to as jumbo frames)
  - This reduces traffic and CPU overhead
- If you use ifconfig to make the changes it does not update ODM so the change does not survive a reboot

	LOGbE Tips
• Low	VCPU entitlement or too few VPs will impact network performance It takes CPU to build those packets
<ul> <li>Con</li> <li>Net</li> <li>abo</li> <li>.</li> </ul>	sider using netperf to test work speed between two LPARs on the same box is limited to the virtual Ethernet Speed which is ut 0.5 to 1.5 Gb/s <u>https://www.ibm.com/developerworks/community/blogs/aixpert/entry/powervm_virtual_ethernet_speed_is_often_confused_with_vios_sea_ive_hea_speed?lang=en</u>
<ul><li>The virt</li><li>But</li></ul>	speed between two LPARs where one is on the SEA and the other is external is the lower of the ual Ethernet speed above or the speed of the physical network all VMs on a server can be sending and receiving at the virtual ethernet speed concurrently
• If 10	OGb network check out Gareth's Webinar
•	http://public.dhe.ibm.com/systems/power/community/aix/PowerVM_webinars/7_10Gbit_Ethernet.wmv Handout at: https://www.ibm.com/developerworks/wikis/download/attachments/153124943/7_PowerVM_10Gbit_Ethernet.pdf?version= 1
22	





# <section-header><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item>



## **Network Commands** • entstat -d or netstat -v (also -m and -I) netpmon • iptrace (traces) and ipreport (formats trace) tcpdump traceroute chdev, Isattr • no ifconfig • ping and netperf or iperf • ftp · Can use ftp to measure network throughput BUT is single threaded • ftp to target • ftp> put "| dd if=/dev/zero bs=32k count=100" /dev/null • Compare to bandwidth (For 1Gbit - 948 Mb/s if simplex and 1470 if duplex ) • 1Gbit = 0.125 GB = 1000 Mb = 100 MB) but that is 100%

27

# netstat -i

netstat -i shows the network interfaces along with input and output packets and errors. It also gives the number of collisions. The Mtu field shows the maximum ip packet size (transfer unit) and should be the same on all systems. In AIX it defaults to 1500.

Both Oerrs (number of output errors since boot) and lerrs (Input errors since boot) should be < 0.025. If Oerrs>0.025 then it is worth increasing the send queue size. Ierrs includes checksum errors and can also be an indicator of a hardware error such as a bad connector or terminator.

The Collis field shows the number of collisions since boot and can be as high as 10%. If it is greater then it is necessary to reorganize the network as the network is obviously overloaded on that segment.

# netstat -i Name Mtu Network Address en6 1500 10.250.134 b740vio1	lpkts lerrs Opkts Oerrs Coll 4510939 0 535626 0 0
# netstat -i         Address           Name         Mtu         Network         Address           en5         1500         link#2         a.aa.69.2b.91.c           en5         1500         10.250.134         b814vio1           lo0         16896         link#1         loopback           lo0         16896         127         loopback           lo0         16896         ::1%1         loopback	Ipkts lerrsOpkts OerrsColl6484659030090610064846590300906100128924401289232001289244012892320012892440128923200
28	



# PCI Adapter transmit Queue Sizes

Table 10. Examples of PCI adapter	transmit queue sizes			
Adapter Type	Feature Code	ODM attribute	Default value	Range
IBM 10/100 Mbps Ethernet PCI Adapter	2968	tx_que_size	8192	16-16384
10/100 Mbps Ethernet Adapter II	4962	tx_que_sz	8192	512-16384
Gigabit Ethernet PCI (SX or TX)	2969, 2975	tx_que_size	8192	512-16384
Gigabit Ethernet PCI (SX or TX)	5700, 5701, 5706, 5707	tx_que_sz	8192	512-16384
10 Gigabit Ethernet PCI-X (LR or SR)	5718, 5719	tx_que_sz	8192	512-16384
ATM 155 (MMF or UTP)	4953, 4957	sw_txq_size	2048	50-16384
ATM 622 (MMF)	2946	sw_txq_size	2048	128-32768
FDDI	2741, 2742, 2743	tx queue size	256	3-2048

For adapters that provide hardware queue limits, changing these values will cause more real memory to be consumed on receives because of the control blocks and buffers associated with them. Therefore, raise these limits only if needed or for larger systems where the increase in memory use is negligible. For the software transmit queue limits, increasing these limits does not increase memory usage. It only allows packets to be queued that were already allocated by the higher layer protocols.

PCI Adapter	Receive	Queue	Sizes

	_				
Table 11	Examples	of PCI	adapter	receive	aueue sizes
10010 111	2/10/10/00	0, , 0,	cicicipitor	1000110	94040 01200

Adapter Type	Feature Code	ODM attribute	Default value	Range
IBM 10/100 Mbps Ethernet PCI Adapter	2968	rx_que_size	256	16, 32 ,64, 128, 26
		rx_buf_pool_size	384	16-2048
10/100 Mbps Ethernet PCI Adapter II	4962	rx_desc_que_sz	512	100-1024
1		rxbuf_pool_sz	1024	512-2048
Gigabit Ethernet PCI (SX or TX)	2969, 2975	rx_queue_size	512	512 (fixed)
Gigabit Ethernet PCI-X (SX or TX)	5700, 5701, 5706, 5707, 5717, 5768, 5271, 5274, 5767, and	rxbuf_pool_sz	2048	512-16384,1
	5281	rxdesc_que_sz	1024	128-3840,128
10 Gigabit PCI-X (SR or LR)	5718, 5719	rxdesc_que_sz	1024	128-1024, by 128
		rxbuf_pool_sz	2048	512-2048
ATM 155 (MMF or UTP)	4953, 4957	rx_buf4k_min	x60	x60-x200 (96-512)
ATM 622 (MMF)	2946	rx_buf4k_min	256 <sup>2</sup>	0-4096
		rx_buf4k_max	0 1	0-14000
FDDI	2741, 2742, 2743	RX_buffer_cnt	42	1-512

31

# txdesc\_que\_sz

Some drivers allow you to tune the size of the transmit ring or the number of transmit descriptors.

The hardware transmit queue controls the maximum number of buffers that can be queued to the adapter for concurrent transmission. One descriptor typically only points to one buffer and a message might be sent in multiple buffers. Many drivers do not allow you to change the parameters.

Adapter type	Feature code	ODM attribute	Default value	Range
Gigabit Ethernet PCI-X, SX or TX	5700, 5701, 5706, 507	txdesc_que_sz	512	128-1024, multiple of 128
32				

# Other Network netstat –v Look for overflows and memory allocation failures Max Packets on S/W Transmit Queue: 884 S/W Transmit Queue Overflow: 9522 "Software Xmit Q overflows" or "packets dropped due to memory allocation failure" Increase adapter xmit queue Use Isattr –El ent? To see setting Look for receive errors or transmit errors dma underruns or overruns mbuf errors

hnimuladov. Clarov	a anto	
ent0 Available 05	i-00 4-Port 10/100/1000 Base-TX P(	CI-Express Adapter (14106803)
bnim: lsattr -El ent0		
chksum_offload yes	Enable hardware transmit and	receive checksum True
flow_ctrl yes	Enable Transmit and Receive Flow (	Control True
jumbo_frames no	Transmit jumbo frames	True
large_send yes	Enable hardware TX TCP resegment	ntation True
rxbuf_pool_sz 2048	Rcv buffer pool, make 2X rxdes	c_que_sz True
rxdesc_que_sz 1024	Rcv descriptor queue size	True
tx_que_sz 8192	Software transmit queue size	True
txdesc_que_sz 512	TX descriptor queue size	True
bnim: Isattr -El en0		
mtu 1500	Maximum IP Packet Size for This Device	ce True
mtu_bypass off	Enable/Disable largesend for virtua	ll Ethernet True
remmtu 576	Maximum IP Packet Size for REMOTE	E Networks True
tcp_nodelay	Enable/Disable TCP_NODELAY Option	n True
thread off	Enable/Disable thread attribute	Тгие

onim: Isattr -El ent5			
na_mode auto High Availability Mode	True		
umbo_frames no Enable Gigabit Ethernet Jumbo Frames	True		
arge_receive no Enable receive TCP segment aggregation	True		
argesend 1 Enable Hardware Transmit TCP Resegmentation	n True		
nthreads 7 Number of SEA threads in Thread mode	True		
ovid 1 PVID to use for the SEA device	True		
ovid_adapter ent4 Default virtual adapter to use for non-VLAN	-tagged packets	True	
queue_size 8192 Queue size for a SEA thread	True		
real_adapter ent0 Physical adapter associated with the SEA	True		
thread 1 Thread mode enabled (1) or disabled (0)	True		
virt_adapters ent4 List of virtual adapters associated with the Si	EA (comma separate	ed) True	
onim: lsattr -El en7			
mtu 1500 Maximum IP Packet Size for This Device T	True		
mtu bypass off Enable/Disable largesend for virtual Etherr	net True		
remmtu 576 Maximum IP Packet Size for REMOTE Netwo	orks True		
tcp nodelay Enable/Disable TCP NODELAY Option	True		

## **Other Network**

- tcp\_nodelayack
  - Disabled by default
  - TCP delays sending Ack packets by up to 200ms, the Ack attaches to a response, and system overhead is minimized
  - Tradeoff if enable this is more traffic versus faster response
  - Reduces latency but increases network traffic
  - The *tcp\_nodelayack* option prompts TCP to send an immediate acknowledgement, rather than the potential 200 ms delay. Sending an immediate acknowledgement might add a little more overhead, but in some cases, greatly improves performance.
  - To set either: chdev -l en0 -a tcp\_nodelay=1
  - OR: no -p -o tcp\_nodelayack=1

## **Other Network**

## • Iparstat 2

• High vcsw (virtual context switch) rates can indicate that your LPAR or VIO server does not have enough entitlement

## ipqmaxlen

37

- netstat -s and look for ipintrq overflows
- ipqmaxlen is the only tunable parameter for the IP layer
- It controls the length of the IP input queue default is 100
- Tradeoff is reduced packet dropping versus CPU availability for other processing

## • Also check errpt – people often forget this

**TCP** Analysis netstat -p tcp tcp: 1629703864 packets sent 684667762 data packets (1336132639 bytes) 117291 data packets (274445260 bytes) retransmitted 955002144 packets received 1791682 completely duplicate packets (2687306247 bytes) 0 discarded due to listener's queue full 4650 retransmit timeouts 0 packets dropped due to memory allocation failure 1. Compare packets sent to packets retransmitted – retransmits should be <5-10% 1. Above is 0.168 2. Compare packets received with completely duplicate packets – duplicates should be <5-10% 1. Above is 2.81 3. In both these cases the problem could be a bottleneck on the receiver or too much network traffic 4. Look for packets discarded because listeners queue is full – could be a buffering issue at the sender 38

# **IP** Stack

ip:

- 955048238 total packets received
- 0 bad header checksums
- 0 fragments received
- 0 fragments dropped (dup or out of space)
- 0 fragments dropped after timeout
- 1. If bad header checksum or fragments dropped due to dup or out of space
  - 1. Network is corrupting packets or device driver receive queues are too small
- 2. If fragments dropped after timeout >0
  - 1. Look at ipfragttl as this means the time to life counter for the ip fragments expired before all the fragments of the datagram arrived. Could be due to busy network or lack of mbufs.
- 3. Review ratio of packets received to fragments received
  - 1. For small MTU if >5-10% packets getting fragmented then someone is passing packets greater than the MTU size

39

# ipqmaxlen

## Default is 100

Only tunable parameter for IP Controls the length of the IP input queue netstat –p ip Look for ipintrq overflows

Default of 100 allows up to 100 packets to be queued up

If increase it there could be an increase in CPU used in the off-level interrupt handler Tradeoff is reduced packet dropping versus CPU availability for other processing

# netstat –v vio

#### SEA Transmit Statistics:

41

Receive Statistics:

Packets: 83329901816 Bytes: 87482716994025 Interrupts: 0 Transmit Errors: 0 Packets Dropped: 0 Packets: 83491933633 Bytes: 87620268594031 Interrupts: 18848013287 Receive Errors: 0 Packets Dropped: 67836309

Bad Packets: 0 Max Packets on S/W Transmit Queue: 374 S/W Transmit Queue Overflow: 0 Current S/W+H/W Transmit Queue Length: 0

Elapsed Time: 0 days 0 hours 0 minutes 0 seconds Broadcast Packets: 1077222 Broad Multicast Packets: 3194318 Multic No Carrier Sense: 0 CRCE DMA Underrun: 0 DMA Lost CTS Errors: 0 Alignr Max Collision Errors: 0 No Re

Broadcast Packets: 1075746 Multicast Packets: 3194313 CRC Errors: 0 DMA Overrun: 0 Alignment Errors: 0 No Resource Errors: 67836309

Virtual I/O Ethernet Adapter (I-lan) Specific Statistics:

Hypervisor Send Failures: 4043136 Receiver Failures: 4043136 Send Errors: 0 Hypervisor Receive Failures: 67836309 "No Resource Errors" can occur when the appropriate amount of memory can not be added quickly to vent buffer space for a workload situation.

You can also see this on LPARs that use virtual Ethernet without an SEA

# Buffers as seen on VIO SEA or Virtual Ethernet

#### # lsattr -El ent5

alt_addr 0x00000000	000 Alternate Ethernet Address	True
chksum_offload yes	Checksum Offload Enable	True
copy_buffs 32	Transmit Copy Buffers	True
copy_bytes 65536	Transmit Copy Buffer Size	True
desired_mapmem 0	I/O memory entitlement rese	erved for device False
max_buf_control 64	Maximum Control Buffers	True
max_buf_huge 64	Maximum Huge Buffers	True
max_buf_large 64	Maximum Large Buffers	True
max_buf_medium 256	Maximum Medium Buffers	True
max_buf_small 2048	Maximum Small Buffers	True
max_buf_tiny 2048	Maximum Tiny Buffers	True
min_buf_control 24	Minimum Control Buffers	True
min_buf_huge 24	Minimum Huge Buffers	True
min_buf_large 24	Minimum Large Buffers	True
min_buf_medium 128	Minimum Medium Buffers	True
min_buf_small 512	Minimum Small Buffers	True
min_buf_tiny 512	Minimum Tiny Buffers	True
poll_uplink no	Enable Uplink Polling	True
poll_uplink_int 1000	Time interval for Uplink Polling	True
trace_debug no	Trace Debug Enable	True
use_alt_addr no	Enable Alternate Ethernet Addr	ess True

Receive Information					
Receive Buffers					
Buffer Type	Tiny	Small	Medium	Large	Huge
Min Buffers	512	512	128	24	24
Max Buffers	2048	2048	256	64	64
Allocated	513	2042	128	24	24
Registered	511	506	128	24	24
History					
Max Allocated	532	2048	128	24	24
Lowest Registered	502	354	128	24	24
(Max Allocated Tepresent (Min Buffers" is number of (Max Buffers" is an absolu	f_small=2048 - f_small=4096	-P -P -P	ny buffers car	n be alloca	ated

netsta	at –p udp
udp:	
	42963 datagrams received
	0 incomplete headers
	0 bad data length fields
	0 bad checksums
	41 dropped due to no socket
	9831 broadcast/multicast datagrams dropped due to no socket
	0 socket buffer overflows
	33091 delivered
	27625 datagrams output
1.	Look for bad checksums (hardware or cable issues)
2	Sockat huffor overflows
<b>∠.</b> .	Could be out of CDU or 1/O bandwidth
1 2 4	<ul> <li>Could be out of CFO of 1/0 bandwidth</li> <li>Could be insufficient UDP transmit or receive sockets, too few nfsd daemons or too small nfs_socketsize or udp_recvspace</li> </ul>





# Performance Tools



DSO (Dynamic System Optimizer)
Not supported on POWER8
<ul> <li>PowerVM and AIX feature         <ul> <li>P7 or P7+</li> <li>AIX v6.1 TL08 SP1 or AIX v7.1 TL02 SP1</li> <li>Cannot be using AMS (Active memory sharing)</li> <li>Chargeable feature via an enablement fileset</li> </ul> </li> </ul>
<ul> <li>Dynamically tunes the allocation of system resources within an LPAR</li> <li>Identifies and optimizes workloads</li> <li>Tries to optimize cache affinity, memory affinity, use of large pages, hardware pre-fetch</li> </ul>
<ul> <li>See chapter 16 of the PowerVM Managing and Monitoring Redbook SG24-7590         <ul> <li><u>http://www.redbooks.ibm.com/redbooks/pdfs/sg247590.pdf</u></li> </ul> </li> </ul>
<ul> <li>Whitepaper at:         <ul> <li><u>http://www-01.ibm.com/common/ssi/cgi-bin/ssialias?infotype=SA&amp;subtype=WH&amp;htmlfid=POW03093USEN</u></li> </ul> </li> </ul>

topas New -L flag for LPAR view nmon nmon analyzer Windows tool so need to co over in ascii mode Opens as an excel spreadsh analyses the data Also look at nmon consolid sar sar -A -o filename 2 30 >/de Creates a snapshot to a file snaps 2 seconds apart Must be post processed on system	ioo, vmo, schedo, vmstat –v lvmo lparstat, mpstat iostat check out Alphaworks for the Graphical LPAR tool eet and then Ganglia - http://ganglia.info Nmonrrd and nmon2web and pGraph Commercial IBM ev/null PM for AIX – in this case 30 Performance Toolbox Tivoli ITM same level of
errpt Check for changes from de	efaults

## Other tools

- filemon
  - filemon -v -o filename -O all
  - sleep 30
  - trcstop
- pstat to check async I/O in 5.3
  - pstat -a | grep aio | wc -l
- perfpmr to build performance info for IBM if reporting a PMR
  - /usr/bin/perfpmr.sh 300

51

## nmon and New Features for V12 Must be running nmon12e or higher • Nmon comes with AIX at 5.3 tl09 or 6.1 tl01 and higher BUT on 5.3 I download the latest version from the web so I get the latest v12 for sure Creates a file in the working directory that ends .nmon This file can be transferred to your PC and interpreted using nmon analyser or other tools Disk Service Times Selecting Particular Disks Time Drift Multiple Page Sizes Timestamps in UTC & no. of digits More Kernel & Hypervisor Stats \* High Priority nmon Virtual I/O Server SEA Partition Mobility (POWER6) WPAR & Application Mobility (AIX6) Dedicated Donating (POWER6) Folded CPU count (SPLPAR) Multiple Shared Pools (POWER6) Fibre Channel stats via entstat 52

	nmon Monitoring
•	nmon -ft –AOPV^dMLW -s 15 -c 120
	<ul> <li>Grabs a 30 minute nmon snapshot</li> <li>A is async IO</li> <li>M is mempages</li> <li>t is top processes</li> <li>L is large pages</li> <li>O is SEA on the VIO</li> <li>P is paging space</li> <li>V is disk volume group</li> <li>d is disk service times</li> <li>^ is fibre adapter stats</li> <li>W is workload manager statistics if you have WLM enabled</li> </ul>
	If you want a 24 hour nmon use:
	nmon -ft –AOPV^dMLW -s 150 -c 576
	May need to enable accounting on the SEA first – this is done on the VIO chdev –dev ent* -attr accounting=enabled
	Can use entstat/seastat or topas/nmon to monitor – this is done on the vios topas –E nmon -O
	VIOS performance advisor also reports on the SEAs
53	











## References Simultaneous Multi-Threading on POWER7 Processors by Mark Funk • http://www.ibm.com/systems/resources/pwrsysperf\_SMT4OnP7.pdf Processor Utilization in AIX by Saravanan Devendran https://www.ibm.com/developerworks/mydeveloperworks/wikis/home?lang=en#/wiki/Power%20Systems/page/Understanding%20CPU%20utilization%20on%20AIX • Gareth Coates – Tricks of the POWER Masters http://public.dhe.ibm.com/systems/power/community/aix/PowerVM webinars/30 Tricks of the Power Masters.pdf Nigel – PowerVM User Group https://www.ibm.com/developerworks/mydeveloperworks/wikis/home?lang=en#/wiki/Power%20Systems/page/PowerVM%20techn ical%20webinar%20series%20on%20Power%20Systems%20Virtualization%20Irom%20IBM%20web • SG24-7940 - PowerVM Virtualization - Introduction and Configuration <u>http://www.redbooks.ibm.com/redbooks/pdfs/sg247940.pdf</u> • SG24-7590 – PowerVM Virtualization – Managing and Monitoring <u>http://www.redbooks.ibm.com/redbooks/pdfs/sg247590.pdf</u> • SG24-8080 – Power Systems Performance Guide – Implementing and Optimizing http://www.redbooks.ibm.com/redbooks/pdfs/sg248080.pdf SG24-8079 – Power 7 and 7+ Optimization and Tuning Guide <u>http://www.redbooks.ibm.com/redbooks/pdfs/sg248079.pdf</u> • Redbook Tip on Maximizing the Value of P7 and P7+ through Tuning and Optimization http://www.redbooks.ibm.com/technotes/tips0956.pdf 59



# Definitions – tcp\_recvspace

tcp\_recvspace specifies the system default socket buffer size for receiving data. This affects the window size used by TCP. Setting the socket buffer size to 16KB (16,384) improves performance over Standard Ethernet and token-ring networks. The default is a value of 4096; however, a value of 16,384 is set automatically by the rc.net file or the rc.bsdnet file (if Berkeley-style configuration is issued).

Lower bandwidth networks, such as Serial Line Internet Protocol (SLIP), or higher bandwidth networks, such as Serial Optical Link, should have different optimum buffer sizes. The optimum buffer size is the product of the media bandwidth and the average round-trip time of a packet. tcp\_recvspace network option can also be set on a per interface basis via the chdev command.

Optimum\_window = bandwidth \* average\_round\_trip\_time

The tcp\_recvspace attribute must specify a socket buffer size less than or equal to the setting of the sb\_max attribute

Settings above 65536 require that rfc1323=1 (default is 0)

61

# Definitions – tcp\_sendspace

tcp\_sendspace Specifies the system default socket buffer size for sending data. This affects the window size used by TCP. Setting the socket buffer size to 16KB (16,384) improves performance over Standard Ethernet and Token-Ring networks. The default is a value of 4096; however, a value of 16,384 is set automatically by the rc.net file or the rc.bsdnet file (if Berkeley-style configuration is issued).

Lower bandwidth networks, such as Serial Line Internet Protocol (SLIP), or higher bandwidth networks, such as Serial Optical Link, should have different optimum buffer sizes. The optimum buffer size is the product of the media bandwidth and the average round-trip time of a packet. tcp\_sendspace network option can also be set on a per interface basis via the chdev command.

Optimum\_window = bandwidth \* average\_round\_trip\_time

The tcp\_sendspace attribute must specify a socket buffer size less than or equal to the setting of the sb\_max attribute

Settings above 65536 require that rfc1323=1 (default is 0)

# Definitions – netstat -m

netstat -m s used to analyze the use of mbufs in order to determine whether these are the bottleneck. The no -a command is used to see what the current values are. Values of interest are thewall, lowclust, lowmbuf and dogticks.

An mbuf is a kernel buffer that uses pinned memory and is used to service network communications. Mbufs come in two sizes - 256 bytes and 4096 bytes (clusters of 256 bytes).

Thewall is the maximum memory that can be taken up for mbufs. Lowmbuf is the minimum number of mbufs to be kept free while lowclust is the minimum number of clusters to be kept free. Mb\_cl\_hiwat is the maximum number of free buffers to be kept in the free buffer pool and should be set to at least twice the value of lowclust to avoid thrashing.

NB by default AIX sets thewall to half of memory which should be plenty. It is now a restricted tunable.

```
# no -a -F | grep thewall
thewall = 1572864
# vmstat 1 1
```

System configuration: lcpu=4 mem=3072MB ent=0.50

net	stat –m – Field meanings
You ca tests to	in use the <b>netstat -Zm</b> command to clear (or zero) the mbuf statistics. This is helpful when running o start with a clean set of statistics. The following fields are provided with the <b>netstat -m</b> command:
Field 1	name Definition
By size	e Shows the size of the buffer.
inuse	Shows the number of buffers of that particular size in use.
calls	Shows the number of calls, or allocation requests, for each sized buffer.
failed	Shows how many allocation requests failed because no buffers were available.
delaye	d Shows how many calls were delayed if that size of buffer was empty and theM_WAIT flag was set by the caller.
free	Shows the number of each size buffer that is on the free list, ready to be allocated.
hiwat	Shows the maximum number of buffers, determined by the system, that can remain on the free list. Any free buffers above this limit are slowly freed back to the system.
freed	Shows the number of buffers that were freed back to the system when the free count when above the hiwat limit.
64	http://www-01.ibm.com/support/knowledgecenter/ssw_aix_71/com.ibm.aix.performance/prftungd_pdf.pdf

# netstat -v - Field meanings

#### **Transmit and Receive Errors**

Number of output/input errors encountered on this device. This field counts unsuccessful transmissions due to hardware/network errors. These unsuccessful transmissions could also slow down the performance of the system.

#### Max Packets on S/W Transmit Queue

Maximum number of outgoing packets ever queued to the software transmit queue. An indication of an inadequate queue size is if the maximal transmits queued equals the current queue size (xmt\_que\_size). This indicates that the queue was full at some point. To check the current size of the queue, use the lsattr -El adapter command (where adapter is, for example, ent0). Because the queue is associated with the device driver and adapter for the interface, use the adapter name, not the interface name. Use the SMIT or the chdev command to change the queue size.

#### S/W Transmit Queue Overflow

Number of outgoing packets that have overflowed the software transmit queue. A value other than zero requires the same actions as would be needed if the Max Packets on S/W Transmit Queue reaches the xmt\_que\_size. The transmit queue size must be increased.

65

http://www-01.ibm.com/support/knowledgecenter/ssw\_aix\_71/com.ibm.aix.performance/prftungd\_pdf.pdf

# netstat -v - Field meanings

#### **Broadcast Packets**

Number of broadcast packets received without any error. If the value for broadcast packets is high, compare it with the total received packets. The received broadcast packets should be less than 20 percent of the total received packets. If it is high, this could be an indication of a high network load; use multicasting. The use of IP multicasting enables a message to be transmitted to a group of hosts, instead of having to address and send the message to each group member individually.

#### DMA Overrun

The DMA Overrun statistic is incremented when the adapter is using DMA to put a packet into system memory and the transfer is not completed. There are system buffers available for the packet to be placed into, but the DMA operation failed to complete. This occurs when the MCA bus is too busy for the adapter to be able to use DMA for the packets. The location of the adapter on the bus is crucial in a heavily loaded system. Typically an adapter in a lower slot number on the bus, by having the higher bus priority, is using so much of the bus that adapters in higher slot numbers are not being served. This is particularly true if the adapters in a lower slot number are ATM adapters.

#### **Max Collision Errors**

66

Number of unsuccessful transmissions due to too many collisions. The number of collisions encountered exceeded the number of retries on the adapter.

http://www-01.ibm.com/support/knowledgecenter/ssw\_aix\_71/com.ibm.aix.performance/prftungd\_pdf.pdf

# netstat -v - Field meanings

### Late Collision Errors

Number of unsuccessful transmissions due to the late collision error.

#### **Timeout Errors**

Number of unsuccessful transmissions due to adapter reported timeout errors.

#### **Single Collision Count**

Number of outgoing packets with single (only one) collision encountered during transmission.

### **Multiple Collision Count**

Number of outgoing packets with multiple (2 - 15) collisions encountered during transmission.

#### **Receive Collision Errors**

Number of incoming packets with collision errors during reception.

#### No mbuf Errors

Number of times that mbufs were not available to the device driver. This usually occurs during receive operations when the driver must obtain memory buffers to process inbound packets. If the mbuf pool for the requested size is empty, the packet will be discarded. Use the netstat -m command to confirm this, and increase the parameter thewall.

67

Γ

http://www-01.ibm.com/support/knowledgecenter/ssw aix 71/com.ibm.aix.performance/prftungd pdf.pdf

Network Spee	d Conversi	on		
	power of 2	bits		= 1
	2^10	1024	=	kilobyte
	2^20	1048576	=	megabyte
	2^30	1073741824	=	gigabyte
	2^40	1.09951E+12	=	terabyte
	2^50	1.1259E+15	=	petabyte
	2^60	1.15292E+18	=	exabyte
	2^70	1.18059E+21	=	zettabyte
	2^80	1.20893E+24	=	yottabyte
	2^90	1.23794E+27	=	lottabyte
To Convert:	See Tab			
bits or Bytes	В			
Kbits or KBytes	К			
Mbits or Mbytes	М			
Gbits or Gbytes	G			
Try cor	verter at: <u>http://www.sp</u>	eedguide.net/conve	rsion.php	2

Converts Gigabit	s or Gigabytes							
1 Kilobyte =	1024	bytes	1 Megabyte =	1048576	bytes	1 gigabyte =	1073741824	bytes
Enter number Gbps:	bytes/sec (Bps)	bytes/min (Bpm)	Kbytes/sec (KBps)	Kbytes/min (KBpm)	Mbytes/sec (MBps)	Mbytes/min (MBpm)	Gbytes/sec (GBps)	Gbytes/min (GBpm)
1	134217728	8053063680	131072	7864320	128	7680	0.125	7.5
	bits/sec (bps)	bits/min (bpm)	Kbits/sec (Kbps)	Kbits/min (Kbpm)	Mbits/sec (Mbps)	Mbits/min (Mbpm)	Gbits/sec (Gbps)	Gbits/min (Gbpm)
	1073741824	64424509440	1048576	62914560	1024	61440	1	60
Enter number GBps:	bytes/sec (Bps)	bytes/min (Bpm)	Kbytes/sec (KBps)	Kbytes/min (KBpm)	Mbytes/sec (MBps)	Mbytes/min (MBpm)	Gbytes/sec (GBps)	Gbytes/min (GBpm)
0.125	134217728	8053063680	131072	7864320	128	7680	0.125	7.5
	bits/sec (bps)	bits/min (bpm)	Kbits/sec (Kbps)	Kbits/min (Kbpm)	Mbits/sec (Mbps)	Mbits/min (Mbpm)	Gbits/sec (Gbps)	Gbits/min (Gbpm)
	1073741824	64424509440	1048576	62914560	1024	61440	1	60

# Definitions – netstat -v

netstat -v is used to look at queues and other information. If Max packets on S/W transmit queue is >0 and is equal to current HW transmit queue length then the send queue size should be increased. If the No mbuf errors is large then the receive queue size needs to be increased.

# netstat -v | grep Queue
Max Packets on S/W Transmit Queue: 0
S/W Transmit Queue Overflow: 0
Current S/W+H/W Transmit Queue Length: 0
Current HW Transmit Queue Length: 0

# netstat -v | grep mbuf No mbuf Errors: 0