



# Dynamically Provisioning Resources to Linux Virtual Servers

Richard Young IBM STG Lab Services 9:30 am Tuesday March 11th 2014 Session 14544

www.SHARE.org



#### **Trademarks & Disclaimer**



The following are trademarks of the International Business Machines Corporation in the United States and/or other countries. For a complete list of IBM Trademarks, see www.ibm.com/legal/copytrade.shtml: AS/400, DB2, e-business logo, ESCON, eServer, FICON, IBM, IBM Logo, iSeries, MVS, OS/390, pSeries, RS/6000, S/390, System Storage, System z9, VM/ESA, VSE/ESA, WebSphere, xSeries, z/OS, zSeries, z/VM.

The following are trademarks or registered trademarks of other companies

Java and all Java-related trademarks and logos are trademarks of Sun Microsystems, Inc., in the United States and other countries. LINUX is a registered trademark of Linux Torvalds in the United States and other countries. UNIX is a registered trademark of The Open Group in the United States and other countries. Microsoft, Windows and Windows NT are registered trademarks of Microsoft Corporation. SET and Secure Electronic Transaction are trademarks of their respective companies.

NOTES: Performance is in Internal Throughput Rate (ITR) ratio based on measurements and projections using standard IBM benchmarks in a controlled environment. The actual throughput that any user will experience will vary depending upon considerations such as the amount of multiprogramming in the user's job stream, the I/O configuration, the storage configuration, and the workload processed. Therefore, no assurance can be given that an individual user will achieve throughput improvements equivalent to the performance ratios stated here.

IBM hardware products are manufactured from new parts, or new and serviceable used parts. Regardless, our warranty terms apply. All customer examples cited or described in this presentation are presented as illustrations of the manner in which some customers have used IBM products and the results they may have achieved. Actual environmental costs and performance characteristics will vary depending on individual customer configurations and conditions. This publication was produced in the United States. IBM may not offer the products, services or features discussed in this document in other countries, and the information may be subject to change without notice. Consult your local IBM business contact for information on the product or services available in your area.

All statements regarding IBM's future direction and intent are subject to change or withdrawal without notice, and represent goals and objectives only. Information about non-IBM products is obtained from the manufacturers of those products or their published announcements. IBM has not tested those products and cannot confirm the performance, compatibility, or any other claims related to non-IBM products. Questions on the capabilities of non-IBM products should be addressed to the suppliers of those products.

Prices subject to change without notice. Contact your IBM representative or Business Partner for the most current pricing in your geography. References in this document to IBM products or services do not imply that IBM intends to make them available in every country. Any proposed use of claims in this presentation outside of the United States must be reviewed by local IBM country counsel prior to such use. The information could include technical inaccuracies or typographical errors. Changes are periodically made to the information herein; these changes will be incorporated in new editions of the publication. IBM may make improvements and/or changes in the product(s) and/or the program(s) described in this publication at any time without notice. Any references in this information to non-IBM Web sites are provided for convenience only and do not in any manner serve as an endorsement of those Web sites. The materials at those Web sites are not part of the materials for this IBM product and use of those Web sites is at your own risk.



# Agenda



The Value of Dynamically Provisioning and Deprovisioning Resources

- Dynamically Adjusting Disk Storage Resources
- Dynamically Adjusting Networking Resources
- Dynamically Adjusting Memory Resources
- Dynamically Adjusting CPU Resources
- 6 Automatically Adjusting Memory and CPU Resources



## Agenda



1 The Value of Dynamically Provisioning and Deprovisioning Resources
2 Dynamically Adjusting Disk Storage Resources
3 Dynamically Adjusting Networking Resources
4 Dynamically Adjusting Memory Resources
5 Dynamically Adjusting CPU Resources
6 Automatically Adjusting Memory and CPU Resources



# **Dynamic Resource Configuration**



- Helps to avoid Linux guest restarts and potential outage/downtime resource allocation changes
- Accommodate unplanned increases in application workload demands or application "enhancements" that consume more than expected resource
- It can allow for more efficient overall hypervisor operation (reduced operational overhead)
- Automated policy based reconfiguration is more responsive than manual adjustments.
- May provide assistance with upgrades by provisioning lower levels of resources both before a virtual server is in production and after it is removed from production.



# Agenda



The Value of Dynamically Provisioning and Deprovisioning Resources
 Dynamically Adjusting Disk Storage Resources
 Dynamically Adjusting Networking Resources
 Dynamically Adjusting Memory Resources
 Dynamically Adjusting CPU Resources
 Automatically Adjusting Memory and CPU Resources



# **Dynamically Provisioning Resources**



		2200								
Device Sub	css -t chan.	3390 DevTyp	pe CU	Туре	Use	PIM	PAM	POM	CHPIDs	
0.0.632a 0.0 0.0.632b 0.0 0.0.632c 0.0 zlnx1:~ # ls CHPID Vary	0.04b6 0.04b7 0.04b8 schp   f Cfg.	3390/0 3390/0 3390/0 nead -0 Type	)c 399 )c 399 )c 399 )c 399 6 Cmg	90/e9 90/e9 90/e9 Share	yes ed	f0 f0 f0	f0 f0 f0	ff 1f ff	16011700 16011700 16011700	00000000 00000000 00000000
0.00 1 0.01 1 0.02 1 0.03 1 zlnx1:~ # ch Vary offline	1 1 1 1 nchp -v e 0.01.	1a 1d 0a 04 0 0.01 . done	2 2 1 1 L	1 1 1 1						
Device Sub	chan.	DevTyp	pe CU	Туре	Use	PIM	PAM	POM	CHPIDs	
0.0.632a 0.0 0.0.632b 0.0 0.0.632c 0.0 zlnx1:~ # ls CHPID Vary	0.04b6 0.04b7 0.04b8 schp   P Cfg.	3390/0 3390/0 3390/0 nead -0 Type	)c 399 )c 399 )c 399 )c 399 6 Cmg	90/e9 90/e9 90/e9 Share	yes ed	f0 f0 f0	f0 f0 f0	ff 1f ff	16011700 16011700 16011700	00000000 00000000 00000000
0.00 1 0.01 0 0.02 1 0.03 1 zlnx1:~ # ch Vary online	1 1 1 1 nchp -v 0.01	1a 1d 0a 04 1 0.01	2 2 1 1 L	1 1 1 1						

- All (non-PCI) IO devices are attached via a defined channel
- In a native LPAR implementation you may need to change the channel (CHPID) state from Linux
- Be aware that lscss does not display the CHPID state
- Use chchp and lschp



- Disk Storage Resource Types
  - ECKD
    - Full Volume
    - z/VM Minidisk
  - SCSI Luns
    - Via z/VM Emulated Device
    - Via Dedicated FCP Device
- All types can be dynamically added
- Can be performed whether in a native LPAR or under z/VM
- General Process
  - Add resource from hypervisor
  - Make new resource available
  - Bring virtual device online
  - Provision as usual







zlnx1:~ # lscss -t 3390 Device Subchan. DevType CU Type Use PIM PAM POM CHPIDs 0.0.632a 0.0.04b6 3390/0c 3990/e9 yes f0 f0 ff 16011700 00000000 0.0.632b 0.0.04b7 3390/0c 3990/e9 f0f0 1f16011700 00000000 zlnx1:~ # cat /proc/cio ignore 0.0.6000-0.0.6329 0.0.632c-0.0.63ff zlnx1:~ # echo free 632c > /proc/cio ignore zlnx1:~ # lscss -t 3390 Device Subchan. DevType CU Type Use PIM PAM POM CHPIDs 0.0.632a 0.0.04b6 3390/0c 3990/e9 yes f0 f0 ff 16011700 00000000 0.0.632b 0.0.04b7 3390/0c 3990/e9 f0 f0 1f16011700 00000000 0.0.632c 0.0.04b8 3390/0c 3990/e9 fO ff f0 16011700 00000000 zlnx1:~ # cat /proc/cio ignore 0.0.6000-0.0.6329 0.0.632d-0.0.63ff zlnx1:~ #

- This example is from a native LPAR implementation
- A cio\_ignore list was used on boot to restrict the available devices
- This list can be dynamically modified to make new devices available
- While a disk example is shown, cio\_ignore applies to all IO devices







```
zlnx1:~ # cat /etc/zipl.conf
# Modified by YaST2. Last modification on Tue Nov 20 21:00:21 CST 2012
[defaultboot]
defaultmenu = menu
###Don't change this comment - YaST2 identifier: Original name: linux###
[SLES11_SP2]
    image = /boot/image-3.0.13-0.27-default
    target = /boot/zipl
    ramdisk = /boot/initrd-3.0.13-0.27-default,0x2000000
    parameters = "root=/dev/disk/by-path/ccw-0.0.632a-part1 TERM=dumb cio_ignore=0.0.6000-0.0.6329,0.0.632c-63ff"
```

- The cio\_ignore list is shown on the kernel parameters line of the zipl.conf
- Be sure to update it with newly (de)provisioned devices as you change the configuration of your system







11 Complete your session evaluations online at www.SHARE.org/AnaheimEval





 201 minidisk still not available to Linux and not shown from a z/VM query virtual

• New storage must be

attached or linked before it can be brought online



RGYLXWS8:/ # vmcp link RG	YLXWS8 201 201	MR	
DASD 0190 3390 P01RES R/C	214 CYI	ON DASD 3F	27 SUBCHANNEL = 000F
DASD 0191 3390 VM1US1 R/C	500 CYI	J ON DASD 3F	10 SUBCHANNEL = $0004$
DASD 0192 3390 LS3F18 R/W	50 CYI	J ON DASD 3F	18 SUBCHANNEL = 0007
DASD 019D 3390 P01RES R/C	292 CYI	J ON DASD 3F	27 SUBCHANNEL = 0010
DASD 019E 3390 P01RES R/C	500 CYI	J ON DASD 3F	27 SUBCHANNEL = 0011
DASD 0200 3390 LS3F52 R/W	10015 CYI	JON DASD 3F	52 SUBCHANNEL = 0006
DASD 0201 3390 LS3F18 R/W	3338 CYI	J ON DASD 3F	18 SUBCHANNEL = 0005
RGYLXWS8:/ # ISCSS Device Subchan. DevTvp	e CU Type Use	PIM PAM POM	CHPIDs
0.0.1000 0.0.0000 1732/0	3 1731/03	80 80 ff	c4000000 00000000
0.0.1001 0.0.0001 1732/0	3 1731/03	80 80 ff	d1000000 00000000
0.0.1002 0.0.0002 1732/0	3 1731/03 yes	80 80 ff	c9000000 00000000
0.0.1003 0.0.0003 1732/0	3 1731/03 yes	80 80 ff	dd000000 00000000
	$\frac{23990}{e9}$	<u>80 80 FF</u>	FF000000 0000000
	c 3990/e9	80 80 ff	££000000 00000000
	c 3990/e9 yes	80 80 ff	ff000000 00000000
	0 3215/00 ves	80 80 ff	ff000000 00000000
0.0.0600 0.0.0009 1732/0	1 1731/01 ves	80 80 ff	
0.0.0601 0.0.000a 1732/0	1 1731/01 ves	80 80 ff	0000000 00000000
0.0.0602 0.0.000b 1732/0	1 1731/01 yes	80 80 ff	0000000 0000000
0.0.000c 0.0.000c 0000/0	0 2540/00	80 80 ff	ff000000 0000000
0.0000 0.0.000d 0000/0	0 2540/00	80 80 ff	ff000000 00000000
0.0.000e 0.0.000e 0000/0	0 1403/00	80 80 ff	ff000000 0000000
0.0.0190 0.0.000f 3390/0	c 3990/e9	80 80 ff	ff000000 00000000
0.0.019d 0.0.0010 3390/0	c 3990/e9	80 80 ff	ff000000 00000000
0.0.019e 0.0.0011 3390/0	c 3990/e9	80 80 ff	ff000000 0000000
RGYLXWS8:/ # chccwdev -e			
Setting device 0.0.0201 d	nline		
Done			



A z/VM "link" makes device available.

Can be performed from Linux via 'vmcp"

Must still be brought online via "chccwdev"



**13** Complete your session evaluations online at www.SHARE.org/AnaheimEval



RGYLXWS8:/ Bus-ID	# lsdasd Status 	Name	Device	Туре	BlkSz	Size	Blocks
0.0.0200 0.0.0201	active active	dasda dasdb	94:0 94:4	ECKD ECKD	4096 4096	7041MB 2347MB	1802700 600840
RGYLXWS8:/ Drive Geome I am going Device Labelli Disk la Disk id Extent Extent Compati Blocksi >> ATTE All data o	<pre># dasdfmt - etry: 3338 C to format t number of de ng device bel entifier start (trk n end (trk no) ble Disk Lay ze NTION! &lt;&lt; f that devic</pre>	<b>b 4096 -f</b> ylinders * he device vice : 0x2 : yes : VOL : 0X0 o) : 0 : 500 out : yes : 409 e will be	/dev/da 15 Head /dev/das 01 1 201 69 6	<b>sdb</b> s = 5 db in	0070 Tr the fol	acks lowing way	:
Type "yes"	to continue	, no will	reave th	e alsk	uncouc	nea.	





RGYLXWS8:/ # fdasd -a /dev/dasdb reading volume label ..: VOL1 reading vtoc .....: ok auto-creating one partition for the whole disk... writing volume label... writing VTOC... rereading partition table... RGYLXWS8:/ #

- Disk storage has been dynamically brought online, formatted, and <u>partitioned</u>
- Put file system on new device
  - mkfs -t ext3 -c /dev/dasdb1
- You could now add to a volume group and LVM to dynamically expand a filesystem without bring the Linux system down
  - pvcreate /dev/dasdb1
  - vgextend VG00 /dev/dasdb1
  - Ivextend -L+1G /dev/VG00/LV01 ; add one more GB to LV
  - ext2online /dev/VG00/LV01
  - resize2fs /dev/VG00/LV01





attach 8a2a to rgylxsp2 as 1000 FCP 8A2A ATTACHED TO RGYLXSP2 1000 Ready; T=0.01/0.01 11:23:47 attach 8b2a to rgylxsp2 as 1001 FCP 8B2A ATTACHED TO RGYLXSP2 1001 Ready; T=0.01/0.01 11:23:55 attach 8c2a to rgylxsp2 as 1002 FCP 8C2A ATTACHED TO RGYLXSP2 1002 Ready; T=0.01/0.01 11:23:59 attach 8d2a to rgylxsp2 as 1003 FCP 8D2A ATTACHED TO RGYLXSP2 1003 Ready; T=0.01/0.01 11:24:04

- Dynamically making the FCP devices available to the guest virtual server
- In an NPIV configuration each device will represent a unique WWPN
- Each WWPN must be zoned to the correct storage resource





SHARE Tethnology - Connections - Results

The new FCP devices are available but must be brought online to Linux



17 Complete your session evaluations online at www.SHARE.org/AnaheimEval



0.0.1000 0.0.0014 1732/03 1731/03	80	80	ff	c4000000	00000000	
0.0.1001 0.0.0015 1732/03 1731/03	80	80	ff	d1000000	00000000	
0.0.1002 0.0.0016 1732/03 1731/03	80	80	ff	c9000000	00000000	
0.0.1003 0.0.0017 1732/03 1731/03	80	80	ff	dd000000	00000000	
rgylxsp2:~ # chccwdev -e 1000-1003						
Setting device 0.0.1000 online						
Done						
Setting device 0.0.1001 online						
Done						
Setting device 0.0.1002 online						
Done						
Setting device 0.0.1003 online						
Done						
rgylxsp2:~ # zfcp_disk_configure 0.0.1002 0x5	50050	763090	08856b	0x4003402	0000000A	1
No configuration file for adapter 0.0.1002						
Configuring FCP disk 500507630908856b:4003402a	100000	0000				
rgylxsp2:~ # lsluns -a						
adapter = 0.0.1002						
port = 0x500507630908856b						
lun = 0x4003402a00000000	/de	ev/sg(	0	Disk	IBM:2107900	)
rgylxsp2:~ #						

- The fcp device range is brought on line with a chccwdev command
- The LUNs are defined via the zfcp\_disk\_configure command
- The FCP device, target storage server WWPN, LUN id, and target state are provided
- Isluns –a confirms the defined LUN





- zfcp\_ping & zfcp\_show are new diagnostic tools when can be helpful when need to configure FCP attached storage
- SLES 11 SP3/RHEL 6.4
- libzfcphbaapi0 package

rgylxsp3:~ # zfcp show 0.0.1000 | more

Interconnect Element Name 0x10000051e364632 Interconnect Element Domain ID 036 Interconnect Element Type Switch Interconnect Element Ports 256 ICE Port 000 Online Attached Port [WWPN/ID] 0x50050763060005b2 / 0x240000 [N\_Port] ICE Port 001 Online Attached Port [WWPN/ID] 0x50050763060045b2 / 0x240100 [N\_Port] ICE Port 002 Online Attached Port [WWPN/ID] 0x500507630310007f / 0x240200 [N\_Port]



19 Complete your session evaluations online at www.SHARE.org/AnaheimEval



rgylxsp2:/etc/udev # cd rule	s.d/	
rgylxsp2:/etc/udev/rules.d #	ls	
40-alsa.rules	52-xpram.rules	79-yast2-drivers.rules
51-dasd-0.0.0200.rules	57-osasnmpd.rules	81-mount.rules
51-packagekit-firmware.rules	59-dasd.rules	81-mptctl.rules
51-qeth-0.0.0600.rules	60-readahead.rules	85-usb_autosuspend_devices.rules
51-geth-0.0.0700.rules	70-kpartx.rules	85-usb_elotouch_wakeup.rules
51-zfcp-0.0.1002.rules	70-persistent-net.rules	99-iwlwifi-led.rules
51-zfcp-0.0.1003.rules	71-multipath.rules	99-pcsc_lite.rules
52-hw_random.rules	77-network.rules	
ravlxsp2./etc/udev/rules d #		

- Confirm the udev entries were made so the definitions are persistent.
- Also make sure your z/VM dedicates exist in the user directory so the devices are available after a restart of the guest virtual server



#### rgylxsp2:~ # fdisk /dev/sda

```
Command (m for help): n
Command action
      extended
      primary partition (1-4)
   q
p1
Partition number (1-4, default 1): 1
First sector (2048-41943039, default 2048):
Using default value 2048
Last sector, +sectors or +size{K,M,G} (2048-41943039, default 41943039):
Using default value 41943039
Command (m for help): w
The partition table has been altered!
Calling ioctl() to re-read partition table.
Syncing disks.
rgylxsp2:~ # pvcreate /dev/sda1
 Physical volume "/dev/sda1" successfully created
rgylxsp2:~ #
```

- At this point you can add the device as you normally would
- Define to the multipather, partition, and place a file system on the device or add to a logical volume



# Agenda



The Value of Dynamically Provisioning and Deprovisioning Resources

- Dynamically Adjusting Disk Storage Resources
- **3** Dynamically Adjusting Networking Resources
  - Dynamically Adjusting Memory Resources
  - Dynamically Adjusting CPU Resources
- 6 Automatically Adjusting Memory and CPU Resources





- Much like dynamically adding disk resources a directory alone does not make the NIC available to Linux.
- Once the NIC is defined there are multiple ways to configure it and some methods vary by Linux distribution.
- Care and planning should be taking when adding additional NIC.
   When adding a new NIC mistakes can cause outages on existing functioning NICs in the same guest.





RGYLXWS8:~ # lsqeth		
Device name	: etnu	
card_type cdev0 cdev1 cdev2 chpid online portname portno state	<pre>: GuestLAN QDIO : 0.0.0600 : 0.0.0601 : 0.0.0602 : 00 : 1 : NET172A : 0 : UP (LAN ONLINE)</pre>	
priority_queueing buffer_count layer2 isolation	: always queue 2 : 64 : 1 : none	
RGYLXWS8:~ # znetconf -c Device IDs	Type Card Type CHPID Drv. Name	State
0.0.0600,0.0.0601,0.0.0602 RGYLXWS8:~ #	1731/01 GuestLAN QDIO 00 qeth eth0	online

This system that has only one NIC and a second NIC will be added





- New NIC added to the zVM user directory
  - Virtual device 700
  - Type QDIO
  - VSWITCH NET172B

dirm for rgylxws8 NICDEF 0700 TYPE QDIO DEV 3 LAN SYSTEM NET172B DVHXMT1191I Your NICDEF request has been sent for processing to DIRMAINT DVHXMT1191I at POKLBS1. Ready; T=0.01/0.02 01:43:35 DVHREQ2288I Your NICDEF request for RGYLXWS8 at \* has been accepted. DVHBIU3450I The source for directory entry RGYLXWS8 has been updated. DVHBIU3424I The next ONLINE will take place immediately. DVHDRC3451I The next ONLINE will take place via delta object directory. DVHRLA3891I Your DSATCTL request has been relayed for processing. DVHBIU3428I Changes made to directory entry RGYLXWS8 have been placed DVHBIU3428I online. DVHREQ2289I Your NICDEF request for RGYLXWS8 at \* has completed; with RC DVHREQ2289I = 0. DVHREQ2288I Your DSATCTL request for DIRMAINT at DVHREQ2288I \* has been accepted. DVHREQ2289I Your DSATCTL request for DIRMAINT at DVHREQ22891 \* has completed; with RC = 0.





- "DEFINE NIC" issued to make the new virtual NIC available to the guest
- Since it was already defined in the user directory it automatically coupled to its virtual switch
- znetconf now shows the new virtual NIC

NIC STACEs the NIC define ni	et unc	ontique de la it	is still	l offli	ne	
RGYLXWS8:~ # vmcp couple 70	0 to sys	stem net172b				
HCPCPL2788E NIC 0700 not co	onnected;	already conne	ected to	VSWIT	CH SYSTEM N	IET172B
Error: non-zero CP response	e for cor	nmand 'COUPLE 7	'00 TO SY	YSTEM	NET172B': #	2788
RGYLXWS8:~ # znetconf -u						
Scanning for network device	es					
Device IDs	Туре	Card Type	CHPID I	Drv.		
0.0.0700,0.0.0701,0.0.0702	1731/01	OSA (QDIO)	01 0	qeth		
RGYLXWS8:~ # znetconf -c						
Device IDs	Туре	Card Type	CHPID I	Drv. N	lame	State
0.0.0600,0.0.0601,0.0.0602	1731/01	GuestLAN QDIO	00 0	qeth e	th0	online
RGYLXWS8:~ #						





- We could use tools such as Yast, netconfig, or redhatconfig-network to configure the interface, but we will use znetconf from s390-tools
- znetconf allows you to configure many different possible attributes of the QDIO device
- Note: znetconf does not create a udev entry
- After executing znetconf the device (not the interface)

```
will be online
RGYLXWS8:~ # znetconf -a 0700 -o layer2=1
Scanning for network devices...
Successfully configured device 0.0.0700 (eth1)
```





- To bring the network interface online you need an ifcfg-ethx script
- If you copy an existing file (such as ifcfg-eth0) you should have only a few changes to make
  - IPADDR, NETMASK, NETWORK
  - \_nm\_name, BROADCAST
- It is highly recommended to put a udev entry in place (/etc/udev/rules.d) so you have a persistent configuration across reboots







- You can activate your new configuration with renetwork restart
- If your new interface configuration breaks your existing network, logon to the 3270 console for the guest and move the ifcfg-ethx script to another directory and reissue your rcnetwork restart command.





## Modifying Attributes of an OSA Network Interface



- Can be performed without restarting the server
- Network interface must be taken offline in many cases
- Don't take offline with chccwdev
- Utilize /sys filesystem interface to take offline/online
- Details documented in the Linux Device Drivers, Features, and Commands manual on DeveloperWorks (See link at end of presentation)



rgylxsp2:~ # lsqeth		
Device name	:	eth0
card_type	:	GuestLAN QDIO
cdev0	:	0.0.0600
cdev1	:	0.0.0601
cdev2	:	0.0.0602
chpid	:	00
online	:	1
portname	:	dontcare
portno	:	0
state	:	UP (LAN ONLINE)
priority_queueing	:	always queue 2
buffer_count	:	64
layer2	:	1
isolation	:	none
Device name	:	eth1
Device name card_type	: 	eth1 GuestLAN QDIO
Device name  card_type cdev0	:	eth1 GuestLAN QDIO 0.0.0700
Device name card_type cdev0 cdev1	:  : : :	eth1 GuestLAN QDIO 0.0.0700 0.0.0701
Device name card_type cdev0 cdev1 cdev2	:  : : : :	eth1 GuestLAN QDIO 0.0.0700 0.0.0701 0.0.0702
Device name card_type cdev0 cdev1 cdev2 chpid	:  : : : : :	eth1 GuestLAN QDIO 0.0.0700 0.0.0701 0.0.0702 01
Device name card_type cdev0 cdev1 cdev2 chpid online	:  : : : : : : :	eth1 GuestLAN QDIO 0.0.0700 0.0.0701 0.0.0702 01 1
Device name card_type cdev0 cdev1 cdev2 chpid online portname	: : : : : : : : : : :	eth1 GuestLAN QDIO 0.0.0700 0.0.0701 0.0.0702 01 1 0
Device name card_type cdev0 cdev1 cdev2 chpid online portname portno	: : : : : : : : : :	eth1 GuestLAN QDIO 0.0.0700 0.0.0701 0.0.0702 01 1 0 0
Device name card_type cdev0 cdev1 cdev2 chpid online portname portno state		eth1 GuestLAN QDIO 0.0.0700 0.0.0701 0.0.0702 01 1 0 0 UP (LAN ONLINE)
Device name card_type cdev0 cdev1 cdev2 chpid online portname portno state priority_queueing		eth1 GuestLAN QDIO 0.0.0700 0.0.0701 0.0.0702 01 1 0 0 UP (LAN ONLINE) always queue 2
Device name card_type cdev0 cdev1 cdev2 chpid online portname portno state priority_queueing buffer_count		eth1 GuestLAN QDIO 0.0.0700 0.0.0701 0.0.0702 01 1 0 0 UP (LAN ONLINE) always queue 2 64
Device name card_type cdev0 cdev1 cdev2 chpid online portname portno state priority_queueing buffer_count layer2		eth1 GuestLAN QDIO 0.0.0700 0.0.0701 0.0.0702 01 1 0 0 UP (LAN ONLINE) always queue 2 64 1



- This system has two network interfaces
- The buffer count on one of them will be increased
- The system will not be brought down
- Only the interface being changed will be stopped



31 Complete your session evaluations online at www.SHARE.org/AnaheimEval



```
rgylxsp2:~ # cd /sys/bus/ccwgroup/drivers/qeth/
rgylxsp2:/qeth # cd 0.0.0700/
rgylxsp2:/0.0.0700 # cat buffer_count
64
rgylxsp2:/0.0.0700 # ifconfig eth1 down
rgylxsp2:/0.0.0700 # chccwdev -d 700
Setting device 0.0.0700 offline
Failed (Invalid argument)
```

- The specific device is found under /sys/bus./ccwgroup/drivers/qeth
- The eth1 interface is stopped
- The attempt to take device 700 offline fails because it must be done via the /sys filesystem





rgylxsp2:/0.0	0.0700 #	echo 128	> buffer_co	unt		Technology - Connections - Pa
-bash: echo:	write e	rror: Open	cation not p	ermitted		
rgylxsp2:/0.0	).0700 #	ls				
blkt	cdev0	chpid	if_name	layer2	performance_	stats
power	S	tate	ungroup			
buffer_count	cdev1	driver	inbuf_size	net	portname	
priority_que	leing s	ubsystem				
card_type	cdev2	hw_trap	isolation	online	portno	
recover	u	event				
rgylxsp2:/0.0	).0700 #	echo 0 >	online			
rgylxsp2:/0.0	).0700 #	echo 128	> buffer_co	unt		
rgylxsp2:/0.0	).0700 #	echo 1 >	online			

- At the top you can see the buffer\_count can not be changed while the device is online
- "echo" a 1 or 0 in the "online" file to control the device state
- This same process can be used to change other attribute, but some such as layer2, may need changes to the udev configuration to be made permanent



rgylxsp2:/0.0.0700 # cat buffer\_count
128

rgylxsp2:/0.0.0700 # lsqeth
Device name

card type	: GuestLAN ODIO
cdev0	: 0.0.0700
cdev1	: 0.0.0701
cdev2	: 0.0.0702
chpid	: 01
online	: 1
portname	: 0
portno	: 0
state	: SOFTSETUP
priority_queueing	: always queue 2
buffer_count	: 128
layer2	: 1
isolation	: none

: eth1



- Only device 0700 shown, device 0600 omitted for readability
- buffer\_count has been changed to the maximum
- At this point the "interface" eth1 just needs to be brought up



## **Modifying Attributes - ethtool**



- "ethtool" can dynamically set or query ethernet interface attributes and statistics
- It can be used to set generic receive offload (GRO), TCP segmentation offload (TSO), checksum operations and other options...
- ethtool –k queries "offload" settings

#### zlnx1:/etc/sysconfig/network # ethtool -k eth0 Offload parameters for eth0: rx-checksumming: on tx-checksumming: off scatter-gather: off tcp-segmentation-offload: off udp-fragmentation-offload: off generic-segmentation-offload: off generic-receive-offload: on large-receive-offload: off rx-vlan-offload: on tx-vlan-offload: on ntuple-filters: off receive-hashing: off



# **Modifying Attributes - ethtool**



- ethtool –K can set offload options
- TCP segmentation offload (TSO) applies only to Layer 3 interfaces
- It can be set dynamically as shown
- GRO can also be controlled but as of kernel 2.6.39 GRO and rx-checksumming are on by default

<pre>zlnx1:/etc/sysconfig/network # ethtool -K eth0 tx on sg on tso on</pre>
<pre>zlnx1:/etc/sysconfig/network # ethtool -k eth0</pre>
Offload parameters for eth0:
rx-checksumming: on
tx-checksumming; on
scatter-gather: on
tcp-segmentation-offload: on
udp-fragmentation-offload: off
generic-segmentation-offload: on
generic-receive-offload: on
large-receive-offload: off
rx-vlan-offload: on
tx-vlan-offload: on
ntuple-filters: off
receive-hashing: off


#### **Modifying Attributes - ethtool**



- Methods by which ethtool changes are made permanent vary by distribution and code levels
- For example, Yast has a place to code ethtool options, <u>but only for</u> <u>ethtool –s</u> and NOT the "-K" offload commands
- ethtool "offload" commands could be made permanent in other ways..
  - /etc/init.d/ scripts

Network Card Setup		
General—Address—Hardware-		-Show visible port identif
Device Name	[Change]	Seconds:
eth0		v 5^
-Kernel Module		
Module NameOptions		
gethâ		
FEthtool options		
Options		
ttele1	[Do als]	
нетрі	[Back]	[Cancel]



#### **Modifying Attributes - ethtool**



- Attempting to code the ethtool offload commands results in an error during interface activation
- Only ethtool –s options should be coded here





#### Agenda



The Value of Dynamically Provisioning and Deprovisioning Resources

- 2 Dynamically Adjusting Disk Storage Resources
- 3 Dynamically Adjusting Networking Resources
- 4 Dynamically Adjusting Memory Resources
- 5 Dynamically Adjusting CPU Resources
- 6 Automatically Adjusting Memory and CPU Resources



- You can dynamically adjust memory for your running Linux system
   making your penguins elastic
- To make memory available you must define it to the LPAR or z/VM before you IPL Linux
- Dynamically addable memory is termed hot plug memory
- Hotplug memory is support was provided by APAR VM64524







 Defining "Reserved" storage to the LPAR will allow you to dynamically add memory to a running Linux server running natively in a partition





SHARE Technology - Connections - Results

===== *	* * * Top of File * * *
===== L	JSER RGYLX0E4 1GYLX0E4 1G 2G G
=====	INCLUDE LINDFLT
=====	CPU 00
=====	CPU 01
=====	CRYPTO APVIRTUAL
=====	IUCV ANY
=====	LOADDEV PORTNAME 5005076306138411
=====	LOADDEV LUN 4011402E0000000

- This z/VM guest has a user directory entry with 1GB of initial memory and 2 GB of maximum memory
- In z/VM, changing the memory size or configuration of a guest causes a storage reset (all storage is cleared)
- If you are running Linux natively in an LPAR without z/VM, you would use reserved storage in the LPAR definition to set aside potential additional memory
- In z/VM, define the memory to be dynamically enabled as "standby" storage



SHARE



- "DEFINE STORAGE 1G STANDBY 1G" issued for this guest
- Issuing a DEFINE STORAGE command causes storage to be cleared
- Anything running at the time of the reset will be immediately terminated without running any shutdown procedures
- This means if you issued this command from a CMS EXEC, CMS is no longer running because storage has been cleared.





SHARE

• Example COMMAND statement in User Directory

USER RGYLX0E1 RGYLX0E1 3G 8G G INCLUDE LINDFLT COMMAND DEFINE STORAGE 2G STANDBY 2G CPU 00 CRYPTO APVIRTUAL IUCV ANY OPTION MAXCONN 128 LINK RGYLXMNT 0191 0191 RR



ICH70001I RGYLX0E1 LAST ACCESS AT 20:23:51 ON THURSDAY, SEPTEMBER 22, 2011 00: NIC 0600 is created; devices 0600-0602 defined 00: z/VM Version 6 Release 1.0, Service Level 1002 (64-bit), 00: built on IBM Virtualization Technology 00: There is no logmsg data 00: FILES: 0001 RDR, NO PRT, NO PUN 00: LOGON AT 20:26:20 EDT THURSDAY 09/22/11 00: STORAGE = 2G MAX = 8G INC = 4M STANDBY = 2G RESERVED = 000: Storage cleared - system reset. z/VM V6.1.0 2010-10-15 11:49

4 Manager Com session evaluations on the at Own. SHARE Rg A Coheim Eval





rgylx0e4:~ # cat	/proc/me	eminfo
MemTotal:	2051920	kB
MemFree:	1877596	kB
Buffers:	10304	kB
Cached:	51160	kB
SwapCached:	0	kB
Active:	29788	kB
Inactive:	54872	kB
Active(anon):	23212	kB
Inactive(anon):	120	kB
Active(file):	6576	kВ
<pre>Inactive(file):</pre>	54752	kВ
Unevictable:	0	kВ
Mlocked:	0	kB
SwapTotal:	0	kB

- After IPLing Linux in this guest, observe via /proc/meminfo that approximately 2GB of memory is available
- The "standby" memory is not reported by /proc/meminfo
- The /sys file system however has an awareness of this "standby" or "hot plug" memory
- With s390-tools, Ismem can be used to report this information and chmem to bring storage elements online or offline





rgylx0e4:~ # lsmem Address Range Removable Size (MB) State Device Core Memory 0 - 63256 online no 0x00000001000000-0x00000006ffffff 1536 online 64-447 ves 0x00000007000000-0x00000007ffffff 256 online 448-511 no 0x00000008000000-0x0000000fffffff 512-1023 2048 offline -----Hotplug Memory Memory device size : 4 MB Memory block size : 256 MB Total online memory : 2048 MB Total offline memory: 2048 MB

- The Ismem command is an easy way to view core and hotplug memory status
- The display looks and works the same whether running under z/VM or running natively



SHARE

rgylx0e4:~ # chmem -e 2g rgylx0e4:~ # lsmem Address Range	Size (MB)	State	Removable	Device
0x000000000000000000000000000000000000	256 1536 256 2048	online online online online	no yes no yes	0-63 64-447 448-511 512-1023
Memory device size : 4 MB Memory block size : 256 MB Total online memory : 4096 MB Total offline memory: 0 MB				

- An additional 2GB of memory now available for use
- The change is temporary, when Linux is restarted, hotplug memory will be offline.
- Remember to make permanent changes for the dynamic resource changes.





rgylx0e4:~ # chmem -d 2g rgylx0e4:~ # lsmem Address Range	Size (MB)	State	Removable	Device
0x000000000000000000000000000000000000	256 1536 256 2048	online online online offline	no yes no -	0-63 64-447 448-511 512-1023
Memory device size : 4 MB Memory block size : 256 MB Total online memory : 2048 MB Total offline memory: 2048 MB	כ			

Storage no longer needed can also be removed to ensure efficient operation



zlnx1:~ Address	# lsmem Range	Size	(MB)	State	Removable	Device
== 0x00000 0x00000 0x00000 0x00000	00000000000-0x000000000000000000000000		256 1024 768 2048	online online online offline	no yes no -	0-1 2-9 10-15 16-31
Memory 1 Total of Total o zlnx1:~ zlnx1:~ Address	block size : 256 MB nline memory : 2048 MB ffline memory: 2048 MB # chmem -e 1024 # lsmem Range	Size	(MB)	State	Removable	Device
== 0x00000 0x00000 0x00000 0x00000 0x00000	00000000000-0x000000000000000000000000		256 1024 768 1024 1024 1024	online online online online offline	no yes no yes -	0-1 2-9 10-15 16-23 24-31
Memory 1 Memory 1 Total o Total o zlnx1:~ chmem: 0 zlnx1:~	device size : 128 MB block size : 256 MB nline memory : 3072 MB ffline memory: 1024 MB # chmem -d 3072 Could only set 2048 MB of memory	y off]	Line.			



 The process and results are the same when running in a native LPAR as shown

Attempts to take more memory offline than possible will result in only the removable memory being taken offline

•



# Dynamically Provisioning Memory – Large Pages

- Large pages can be added permanently via hugepages=<npages> in the kernel parameter line of zipl.conf
- Huge page information can be queried via /proc/meminfo

HugePages_Total:	0
HugePages_Free:	0
HugePages_Rsvd:	0
HugePages_Surp:	0
Hugepagesize:	1024 kB

- Also queried via /proc/sys/vm/nr\_hugepages
- Can be set dynamically via echo xxx > /proc/sys/vm/nr\_hugepages
- Hotplug memory allocated as moveable and can only be used by movable resources.
- By default Large Pages are not allocated as movable resource but can be made to allocate from movable

51 poto were of yliwith Share.org/AnaheimEval



# Dynamically Provisioning Memory – Large Pages

VmallocChunk: 134150024 kB
HugePages_Total: 0
HugePages_Free: 0
HugePages_Rsvd: 0
HugePages_Surp: 0
Hugepagesize: 1024 kB
<pre>zlnx1:~ # cat /proc/sys/vm/nr_hugepages</pre>
0
<pre>zlnx1:~ # echo 1500 &gt; /proc/sys/vm/nr_hugepages</pre>
<pre>zlnx1:~ # cat /proc/sys/vm/nr_hugepages</pre>
1500
zlnx1:~ #

VmallocChunk:	134150024	kB
HugePages_Total:	1500	
HugePages_Free:	1500	
HugePages_Rsvd:	0	
HugePages_Surp:	0	
Hugepagesize:	1024	kВ
zlnx1•~ #		



echoalogy - Connections - Recult

 Allocate your large pages as soon as possible to avoid fragmentation issues





- Storage-class memory is a class of data storage devices that combines properties of both storage and memory.
- Storage-class memory(SCM) is represented as a block device. Therefore it could be utilized as swap device or as part of a filesystem, even in a LVM
- A SCM device can be partitioned in to 7 unique partitions
- Requires kernel 3.6 or 3.7 and above, defined to the LPAR, and the scm\_block module loaded
- Use Isblk and Isscm to view resource
- Can use storage commands such as fdisk, mkswap, mkfs, mount with the devices

<pre>zlnx:~ # lsscm SCM Increment</pre>	Size	Name	Rank	D_state	0_state	Pers	ResID	
000000000000000000000000000000000000000	16384MB 16384MB	scma scmb	 1 1	2 2 2	 1 1	 2 2	 1 1	-





#### Hardware Management Console

(a) (b) (b) (c) (c) (c) (c) (c) (c) (c) (c) (c) (c	Tasks Index		Manage Flash Allocation - SCZP401					i
<ul> <li>Welcome</li> <li>Systems Management</li> <li>Ensemble Management</li> <li>TISO Ensemble</li> <li>Members</li> <li>SCZP301</li> <li>SCZP401</li> <li>Workloads</li> </ul>	Image: Second	Summary Allocated: 656 GB Storage increment: 16 GB Available: 2192 GB Rebuild complete: 0 % Uninitialized: 0 GB Unavailable: 0 GB Total: 2848 GB						
📕 HMC Management		Partitic	ons					
Service Management			Select Act	ion	•			
Tasks Index		Select	Partition Name	Status	IOCDS	Allocated (GB)	Maximum (Gl	B)
		۲	A01	Active	A0,A1,A2,A3	16	48	<b>^</b>
		0	A03	Active	A0,A1,A2,A3	32	32	
		0	A08	Active	A0,A1,A2,A3	32	32	=
		0	A09	Active	A0,A1,A2,A3	32	32	
		0	A12	Active	A0,A1,A2,A3	16	32	
		0	A15	Active	A0,A1,A2,A3	16	32	
		0	A19	Active	A0,A1,A2,A3	192	208	
		0	A21	Active	A0,A1,A2,A3	32	64	
SCM storage add	to I DAD	0	A22	Active	A0,A1,A2,A3	32	64	
		0	A23	Active	A0,A1,A2,A3	32	64	-
dynamically		Refres OK	h Apply Cance	Hel	p			





- Load Module modprobe scm\_block
- List block devices
- Make node in /dev filesystem



• Use as needed...





56 Complete your session evaluations online at www.SHARE.org/AnaheimEval



#### **Dynamic Memory - Considerations**



- To add and remove memory takes some small advanced planning. Develop a standard policy around how you will handle memory needs.
- Memory can be added or removed whether you are running under z/VM or in a native LPAR
- zVM User Directory COMMAND statement provides an effective way to issue the DEFINE STORAGE command in an non-disruptive manner.
- Remember not all memory sections will be removable, and the removable state can change over time
- Hot plugged memory is NOT currently managed by cpupluged memory management (cmm)



#### **Summary of Memory Hotplug**



- Basic memory hotplug requirements:
  - ✓ VM64524 support
  - DEFINE STORAGE STANDBY issued before Linux is IPLed
  - For native LPAR, RESERVED STORAGE must be defined before the LPAR is activated
  - ✓ SLES 11 / RHEL 6 provide support in Linux
- Suspend/Resume restriction: The Linux instance must not have used any hotplug memory since it was last booted. (Has worked if freed in advance)
- You may not be able to disable hotplug memory that has been enabled
- Can be very helpful when exact future memory need is unknown, without over allocating online memory from the start.
- After a Linux reboot core memory is made available again and hotplug memory is freed



#### Agenda



The Value of Dynamically Provisioning and Deprovisioning Resources

- Dynamically Adjusting Disk Storage Resources
- 3 Dynamically Adjusting Networking Resources
  - Dynamically Adjusting Memory Resources

5 Dynamically Adjusting CPU Resources

6 Automatically Adjusting Memory and CPU Resources



# **Dynamically Provisioning CPU Resources**



- Multithreaded application or multiple applications in a single virtual server could potentially benefit from additional virtual CPs
- Conversely too many virtual CPs could be counter productive from a performance perspective. If you are over provisioned you can dynamically remove processor resource.
- Adding or removing virtual CPs could impact monitoring applications or middleware that might query the number of processors on startup (ie the Java Virtual Machine)
- z/VM "DEFINE CPU" is a Class G command
- (R.O.T.) Don't add unnecessary virtual CPs and never more virtual CPs than logical processors available.

<sup>60</sup> Complete your session evaluations online at your State or all CPS does not add physical



=====	USER RGYLX0E4 1GYLX0E4 1G 2G G
=====	INCLUDE LINDFLT
=====	CPU 00
=====	CPU 01
=====	CRYPTO APVIRTUAL
=====	IUCV ANY
=====	LOADDEV PORTNAME 5005076306138411
=====	LOADDEV LUN 4011402E00000000
=====	MACHINE ESA 4
=====	OPTION APPLMON MAXCONN 128

- The directory entry shows two initial virtual CPs
- The maximum potential virtual CPs shown is four
- z/VM does not make the additional potential virtual CPs available for Linux to enable on its own
- The additional potential virtual CPs must first be **defined** in the z/VM guest before dynamically enabling on Linux





rgylx0e4:~ # vmcp q v STORAGE = 1G
 XSTORE = none
CPU 00 ID FF12EBBE20978000 (BASE) CP CPUAFF ON
CPU 01 ID FF12EBBE20978000 CP CPUAFF ON
AP 51 CEX2A Queue 08 shared
CONS 0009 DISCONNECTED TERM START
0009 CL T NOCONT NOHOLD COPY 001 READY FORM STANDARD
0009 TO RGYLX0E4 RDR DIST RGYLX0E4 FLASHC 000 DEST OFF
0009 FLASH CHAR MDFY 0 FCB LPP OFF
0009 3215 NOEOF OPEN 0013 NOKEEP NOMSG NONAME
0009 SUBCHANNEL = 000A

- The current z/VM guests virtual resources are displayed from within Linux
- The two initial and active virtual CPs are shown
- Notice there is no information displayed about the potential additional virtual CPs





rgylx0e4:~ Linux 2.6.3	# mpsta 2.29-0.	t -A 3-defaul	t (rgylx	0e4) 04	4/01/11	_ <sup>s</sup>	390x_			
13:19:24	CPU	%usr	%nice	%sys %	iowait	%irq	%soft	%steal	%guest	%idle
13:19:24	0	1.62	0.00	0.67	0.29	0.00	0.02	0.03	0.00	97.33
13:19:24	1	1.25	0.00	0.64	0.30	0.00	0.02	0.08	0.00	97.70
13:19:24	CPU	intr/s								
13:19:24	all	0.00								
13:19:24	0	0.00								
13:19:24	1	0.00								

- Note the mpstat output from before defining the additional virtual CPs
- Observe the even distribution of idle time and usage





```
rgylx0e4:/sys/devices/system/cpu # ls
cpu0 cpu1 dispatching kernel_max offline online perf_events possible present
rgylx0e4:/sys/devices/system/cpu # cat kernel_max
63
rgylx0e4:/sys/devices/system/cpu # cat online
0-1
rgylx0e4:/sys/devices/system/cpu # cat offline
2-63
rgylx0e4:/sys/devices/system/cpu # cat possible
0-63
rgylx0e4:/sys/devices/system/cpu # cat present
0-1
rgylx0e4:/sys/devices/system/cpu # cat sched_mc_power_savings
0
rgylx0e4:/sys/devices/system/cpu #
```

- The Linux sysfs file system can access information about the two active virtual CPs
- The kernel has a maximum potential of 64 processors
- No information about the two potential additional virtual CPs is shown yet





rgylx0e4:/sys/devices/system/cpu # modprobe vmcp rgylx0e4:/sys/devices/system/cpu # vmcp define CPU 03 type cp CPU 03 defined rgylx0e4:/sys/devices/system/cpu # vmcp define CPU 02 type cp CPU 02 defined rgylx0e4:/sys/devices/system/cpu # ls cpu0 cpu1 dispatching kernel\_max offline online perf\_events possible rgylx0e4:/sys/devices/system/cpu #

- Using the vmcp command we pass the zVM CP DEFINE CPU commands on to our z/VM guest.
- Remember this is a class G guest enabling the additional resources previously defined in the user directory
- After defining the additional virtual CPs in z/VM we still do not see them in the Linux /sys filesystem.





rgylx0e4:/sys/devices/system/cpu # ls cpu0 cpu1 dispatching kernel max offline online perf events possible present rescan rgylx0e4:/sys/devices/system/cpu # vmcp g v STORAGE = 1GXSTORE = noneCPU 00 ID FF12EBBE20978000 (BASE) CP CPUAFF ON FF12EBBE20978000 CP CPU 01 TD CPUAFF ON CPU 03 ID FF12EBBE20978000 STOPPED CP CPUAFF ON CPU 02 ID FF12EBBE20978000 STOPPED CP CPUAFF ON AP 51 CEX2A Oueue 08 shared CONS 0009 DISCONNECTED TERM START 0009 CL T NOCONT NOHOLD COPY 001 READY FORM STANDARD 0009 TO RGYLX0E4 RDR DIST RGYLX0E4 000 DEST OFF FLASHC 0009 FLASH CHAR MDFY 0 FCB LPP OFF 0009 3215 NOEOF OPEN 0013 NOKEEP NOMSG NONAME 0009 SUBCHANNEL = 000ARDR 000C CL \* NOCONT NOHOLD EOF READY 000C 2540 CLOSED NOKEEP NORESCAN SUBCHANNEL = 000E

- By using the z/VM QUERY VIRTUAL command we can see the additional virtual CPs have been defined to the guest
- The new virtual CPs are in a "stopped" state





rgylx0e4:/sys/devices/system/cpu # mpstat -A Linux 2.6.32.29-0.3-default (rgylx0e4) 04/01/11 _s390x_												
13:23:58	CPU	%usr	%nice	%sys %	iowait	%irq	%soft	%steal	%guest	%idle		
13:23:58	all	0.47	0.00	0.23	0.10	0.00	0.01	0.02	0.00	99.16		
13:23:58	0	0.54	0.00	0.24	0.10	0.00	0.01	0.01	0.00	99.10		
13:23:58	1	0.41	0.00	0.23	0.10	0.00	0.01	0.03	0.00	99.23		
rgylx0e4:/sy	rgylx0e4:/sys/devices/system/cpu # ls											
cpu0 cpu1	dispate	hing ke	rnel_max	offline	online	perf_ev	vents po	ossible	present	rescan	sched_mc_p	
rgylx0e4:/sy	ys/devic	es/syste	m/cpu #	echo 1 >	rescan							
rgylx0e4:/sys/devices/system/cpu # ls												
cpu0 cpu1 cpu2 cpu3 dispatching kernel max offline online perf events possible present rescan												
rgylx0e4:/sy	ys/devic	es/syste	m/cpu #									

- mpstat is only reporting two CPUs
- The **rescan** operation is used to search for new CPUs in the guest.
- After rescan, additional / sysfs entries exist





rgylx0e4:/sys/devices/system/cpu # mpstat -A Linux 2.6.32.29-0.3-default (rgylx0e4) 04/01/11 _s390x_											
13:24:41	CPU	%usr	%nice	%sys	%iowait	%irq	%soft	%steal	%guest	%idle	
13:24:41	all	0.43	0.00	0.21	0.09	0.00	0.01	0.02	0.00	99.23	
13:24:41	0	0.49	0.00	0.22	0.09	0.00	0.01	0.01	0.00	99.18	
13:24:41	1	0.37	0.00	0.21	0.09	0.00	0.01	0.02	0.00	99.29	
13:24:41	2	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
13:24:41	3	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	

- mpstat reports 0% use and 0% idle for the new CPUs. This is because they are stopped and offline
- The new CPUs must still be brought online to Linux





```
rgylx0e4:/sys/devices/system/cpu/cpu2 # echo 1 > online
rgylx0e4:/sys/devices/system/cpu/cpu2 # ls
address capability configure crash_notes idle_count idle_time_us online polarization
rgylx0e4:/sys/devices/system/cpu/cpu2 # cat online
1
rgylx0e4:/sys/devices/system/cpu/cpu2 # echo 1 > ../cpu3/online
```

 Bring the new CPUs online to Linux by echoing 1 in to the "online" file for the given CPU





rgylx0e4:/s	sys/devi	ces/syst	tem/cpu # lt (rawly	mpstat	-A		~200 <del>v</del>			
LINUX 2.0.3	52.29-0.	S-derau.	тс (туутх	-	5590X_					
13:26:36	CPU	%usr	%nice	%sys	%iowait	%irq	%soft	%steal	%guest	%idle
13:26:36	all	0.33	0.00	0.17	0.07	0.00	0.01	0.02	0.00	99.41
13:26:36	0	0.39	0.00	0.18	0.07	0.00	0.01	0.01	0.00	99.33
13:26:36	1	0.30	0.00	0.17	0.07	0.00	0.01	0.02	0.00	99.43
13:26:36	2	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	100.00
13:26:36	3	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	100.00

- On a idle system, the new CPUs momentarily show 100% idle after being brought online
- Once a little bit of workload hits the system, this quickly changes





<pre>rgylx0e4:/sys/devices/system/cpu # ls cpu0 cpu1 cpu2 cpu3 dispatching kernel_max offline online perf_events possible rgylx0e4:/sys/devices/system/cpu # echo 0 &gt; cpu1/online rgylx0e4:/sys/devices/system/cpu # echo 0 &gt; cpu3/online rgylx0e4:/sys/devices/system/cpu # mpstat -A</pre>												
Linux 2.6.3	Linux 2.6.32.29-0.3-default (rgylx0e4) 04/01/11 _s390x_											
13:27:53 13:27:53 13:27:53 13:27:53 13:27:53 13:27:53	CPU all 0 1 2 3	%usr 0.27 0.35 0.00 0.00 0.00	%nice 0.00 0.00 0.00 0.00 0.00	%sys 0.14 0.16 0.00 0.00 0.00	%iowait 0.06 0.06 0.00 0.00 0.00	%irq 0.00 0.00 0.00 0.00 0.00	%soft 0.01 0.01 0.00 0.00 0.00	<pre>%steal     0.01     0.01     0.00     0.00     0.00     0.00</pre>	%guest 0.00 0.00 0.00 0.00 0.00	%idle 99.52 99.40 0.00 100.00 0.00		

•You can take dynamically added CPUs offline again

•You can take offline CPUs that were initially online as well





• The latest levels of s390-tools has two new commands

• lscpu													
- ,	[roo	t@rgy	ylxr64 ~	~]# l:	scpu								
• chcpu	Arch	itect	ure:		s390z	s390x							
	CPU	op-mo	ode(s):		32-bi	32-bit, 64-bit							
	Byte	Orde	er:		Big H	Indian							
	CPU (	s):			2								
	On-1	ine (	CPU(s) ]	list:	0,1								
	Thre	ad(s)	per co	ore:	1								
	Core	(s) p	per socl	ket:	1	1							
	Sock	et(s)	per bo	ook:	1	1							
	Book	(s):	-		2	2							
	Vend	or II	):		IBM/S	IBM/S390							
	Bogo	MIPS:			9765.	.00							
	Hype	rviso	or:		z/VM	6.2.0							
	Hype	rviso	or vendo	or:	IBM								
	Virt	ualiz	zation t	cype:	full	full							
	Disp	atchi	ing mode	∋:	horiz	horizontal							
	[roo	tergy	/lxr64 ~	~]# 1:	scpu -ae	;							
	CPU	BOOK	SOCKET	CORE	ONLINE	CONFIGURED	POLARIZATION	ADDRESS					
	0	0	0	0	yes	yes	horizontal	0					
	1	1	1	1	yes	yes	horizontal	1					


### **Dynamically Managing Virtual CPs**







### **Dynamically Managing Virtual CPs**



[root@rgylxr64 ~]# lscpu -ae CPU BOOK SOCKET CORE ONLINE CONFIGURED POLARIZATION ADDRESS horizontal 0 0 0 0 yes yes horizontal 1 1 yes yes 1 1 [root@rgylxr64 ~]# chcpu -p vertical chcpu: Failed to set vertical dispatch mode: Operation not supported [root@rgylxr64 ~]#

- Polarization Configures dispatching in a Hiperdispatch like manner
- Function applies to native LPAR deployments
- Horizontal (default) spread even across all logical process
- Vertical dispatch across as few as possible



### Agenda



Value of Dynamic Resource Configuration
Dynamically Adjusting Disk Storage Resources
Dynamically Adjusting Networking Resources
Dynamically Adjusting Memory Resources
Dynamically Adjusting CPU Resources

### 6 Automatically Adjusting Memory and CPU Resources



### What is cpuplugd and why is it important



- Manually adjusting the quantity of CPU and memory configured to virtual guests is not the most effective approach, especially when managing thousands of virtual servers.
  - Rules based Linux automation for adding and removing memory and processor resources
  - The daemon checks the system at user configurable intervals
- ✓ You must configure the rules for it to operate

 $\checkmark$ 

- ✓ You must activate the cpuplug daemon to use it, by default it is inactive
  - New capabilities have recently been added to cpuplugd with s390-tools 1.15 (RHEL 6.2 & SLES 11 SP2)



### cpuplugd – Planning



- The default rules are NOT recommendations, they are syntax examples.
- You should customize the configuration to fit your environment. Each virtual server may have different needs based on workload, middleware, and other factors.
- cpuplugd -V -f -c /etc/sysconfig/cpuplugd This invokes cpuplugd in the foreground with verbose messaging to help you understand its operation. It is highly recommend you use this to understand how cpuplugd is functioning
  - Send to logfile: cpuplugd -c <config file> -f -V>&<logname> &
- When building rules for cpuplugd, it is important to understand what state you will be in after you execute a "plug" or "unplug" operation when writing the rules.
- Suggested Reading:
  - May 2012 Paper ZSW03228 "Using the Linux cpuplugd Daemon to manage CPU and memory resources from z/VM Linux guests"



### cpuplugd – Planning



- Remember some middleware such as DB2 and Oracle have memory managers and resource optimizing code of their own
- The purpose and operation of these are different than cpuplugd
- With that said any CPUs or memory brought online dynamically would immediately be available for use





### cpuplugd – Rule Considerations



- Ensure you can grow CPU capacity to what the application requires to perform well (don't artificially limit). Use other mechanisms to throttle MIP usage based on shares/priorities.
- Rules based on the last couple of sample intervals are more responsive than ones based on averages over minutes. Slower responses to change can mean lower throughput for your applications
- Keep in mind you can only add/remove a full virtual CP of capacity.
- Avoid rules that plug and immediately unplug CPUs continuously
  - Plug = idle < 50
  - Unplug = idle > 50
- This means at times you might have > 1 virtual CPs of idle capacity as an acceptable state.



### cpuplugd - What if I run with default rules?



- CPU\_MIN= 1 and CPU\_MAX= 0 (maximum available)
- UPDATE= 5
- HOTPLUG="(loadavg > onumcpus + 0.75) & (idle < 10.0)"
- HOTUNPLUG="(loadavg < onumcpus 0.25) | (idle > 50)"
- Basic variables can be defined as:
  - loadavg: The load average over the past minute
  - onumcpus: The number of cpus which are online now
  - runable\_proc: The current quantity of runable processes
  - idle: The current idle percentage
- Unplug at 51% idle? After unplug, what is my cpu busy?
- Plug only at 91% busy? What if my runable processes are growing high?



### cpuplugd – What the variables represent



- idle: Current idle Where 1 idle processor = 100 and 4 idle processors = 400 (/proc/stat 4<sup>th</sup> value). Idle does NOT stop at 100!
- loadavg: The current load average The first /proc/loadavg value. The average number of runnable process. Not average CPU utilization! One looping process on a system would cause this to approach 1.0 Five looping processes on a single CPU system would cause this to approach 5.0
- **onumcpus:** The actual number of cpus which are online (Via: /sys/devices/system/cpu/cpu%d/online)
- runable\_proc: The current quantity of runnable processes (The 4<sup>th</sup> /proc/loadavg value)



### cpuplugd – Variables and rule capabilities for CP4

- New predefined keywords
  - user the current CPU user percentage
  - nice the current CPU nice percentage
  - system the current CPU system percentage
  - idle the current CPU idle percentage
  - iowait the current CPU iowait percentage
  - irq the current CPU irq percentage
  - softirg the current CPU softirg percentage
  - steal the current CPU steal percentage
  - guest the current CPU guest percentage
  - guest\_nice the current CPU guest\_nice percentage
  - cpustat.<name> data from /proc/stat and /proc/loadavg
  - time floating point timestamp in "seconds.microseconds" since Unix Epoch
- Historical function available and extremely useful
  - 0 is current interval
  - cpustat.idle[0] .... cpustat.idle[99]
- User Defined Variables Now Supported (See examples next slide)





### User defined variables example for CPU



- user\_0="(cpustat.user[0] cpustat.user[1])"
- nice\_0="(cpustat.nice[0] cpustat.nice[1])"
- system\_0="(cpustat.system[0] cpustat.system[1])"
- user\_2="(cpustat.user[2] cpustat.user[3])"
- nice\_2="(cpustat.nice[2] cpustat.nice[3])"
- system\_2="(cpustat.system[2] cpustat.system[3])"
- CP\_Active0="(user\_0 + nice\_0 + system\_0)/ (cpustat.total\_ticks[0] cpustat.total\_ticks[1])"
- CP\_Active2="(user\_2 + nice\_2 + system\_2)/ (cpustat.total\_ticks[2] cpustat.total\_ticks[3])"
- CP\_ActiveAVG="(CP\_Active0+CP\_Active2) / 2"
- idle\_0="(cpustat.idle[0] cpustat.idle[1])"
- iowait\_0="(cpustat.iowait[0] cpustat.iowait[1])"
- idle\_2="(cpustat.idle[2] cpustat.idle[3])"
- iowait\_2="(cpustat.iowait[2] cpustat.iowait[3])"
- CP\_idle0="(idle\_0 + iowait\_0)/ (cpustat.total\_ticks[0] cpustat.total\_ticks[1])"
- CP\_idle2="(idle\_2 + iowait\_2)/ (cpustat.total\_ticks[2] cpustat.total\_ticks[3])"
- CP\_idleAVG="(CP\_idle0 + CP\_idle2) / 2"





### cpuplugd memory management features



# Why do we care about dynamic adjustments to memory?



- Too little free memory
  - Can't start new programs/processes
  - Swapping degrades performance
  - OOM killer kills your middleware server process to "save" the system
- Too much free memory
  - Excessive use of page cache, that may get paged out at hypervisor layer
  - Cause stealing of memory pages from other guest with legitimate need
  - · More memory stress and paging at the hypervisor layer
  - Degrades overall performance
- Unlike physical only environments, over allocating memory in virtual environments can be very counter productive. This is not just true for z/VM.
- Manual adjustments don't happen fast enough



### Automated adjustments of memory



- Problems stemming from over/undersized memory allocations of guests in virtualized environments are not unique to Linux on System z
- Even the most accurate sizing is irrelevant as soon as the requirements change
- The cpuplug daemon determines how much memory to add or remove based upon the rules you put in place
- It is based on the same configurable interval you set for CPU rules
- The memory increment added or removed is configurable (and you should)
- Separate plug and unplug rules are used for memory management
- There are NO default memory plug and unplug rules
- If you start cpuplugd without any configuration changes it will manage CPUs but NOT memory.
- Be sure to have the following z/VM PTFs on:
  - APAR VM65060 **REQUIRED**!
    - 540 <u>UM33537</u>
    - 620 <u>UM33539</u>



### Automated adjustments of memory



- cpuplugd uses CMM (or cmm1) to return unused pages of memory to the hypervisor via a diagnose call. This memory must be actually free and not used as cache. The decision as to how many pages to return is controlled by the rules you write for cpuplugd
- The cmm module must be loaded
- Don't mix cpuplugd and VMRM management of CMM.
- CMMA (or cmm2) is an alterative mechanism to return pages of memory to the hypervisor by checking a bit on the page. It will only operate on free pages and z/VM has to perform a scan for memory before the pages are actually reclaimed from a guest.
- Linux pagecache can be the large consumer of free memory. If you guest idles while holding page cache z/VM could page this memory out causing long delays when activity resumes. (Example: slow ssh logins first thing in the morning)
- Understand that manually freeing pagecache alone, does not return the formerly used pages to z/VM. One of the two above CMM mechanisms must be used to return the pages.
- If the pages are not returned to z/VM, they will likely be paged out because they are idle as Linux is not using them. The next time Linux allocates something to that page, it would have to be paged in causing unnecessary delays



### Linux memory management at a high level



- Understanding Linux memory management effects how you might write your plugd rules
- Application requests for memory are managed as follows:
  - > With sufficient free pages, the request is fulfilled immediately
  - If that causes the amount of free memory to fall below a high water mark, an **asynchronous** page scan by kswapd is triggered in the background.
  - If serving the request would cause the amount of free memory to fall below a low water mark, a so called direct scan is triggered, and the application waits until this scan provides the required pages.
  - The system may decide to mark anonymous pages (pages that are not related with files on disks) for swapping and initiate that these pages be written to swap asynchronously.
- The async page scan is in an early indicator of a memory shortage
- Direct scans are more costly in terms of application performance
- Writing rules based on the scans can be more responsive than waiting until some paging activity occurs.



### Automated adjustments of memory

- Basic variables for writing memory plug and unplug rules
  - > apcr: the amount of page cache reads listed in vmstat bi/bo
  - **Freemem:** the amount of free memory (in megabyte)
  - > swaprate the number of swapin & swapout operations
  - cpustat.<name> from /proc/stat and /proc/loadavg
  - > meminfo.<name> any value from /proc/meminfo
  - vmstat.<name> any value from /proc/vmstat
  - > time floating point timestamp in "seconds.microseconds"
- CMM pool size and increment
  - > **CMM\_MIN** min size of static page pool (default 0)
  - CMM\_MAX max size of static page pool
    - default 512MB
  - > CMM\_INC amount for memunplug only (previously for plug and unplug)
    - > 10% of free memory + cache, in pages
  - CMM\_DEC amount for memplug operation
    - > default 10% of total memory in pages
- With heavier IO rates you may want to allow the system to utilize more memory to help improve performance. This memory would get utilized by pagecache.
- Looking at "cache" for free memory might be skewed if you have a lot of shared memory (databases or java for example)



Anaheim

### User defined variable example for memory



- The page scan rate can be calculated as the sum of:
  - vmstat.pgscan\_kswapd\_dma
  - vmstat.pgscan\_kswapd\_normal
  - vmstat.pgscan\_kswapd\_movable
    - pgscan\_k="vmstat.pgscan\_kswapd\_dma[0] + vmstat.pgscan\_kswapd\_normal[0] + vmstat.pgscan\_kswapd\_movable[0]"
- The direct page scan rate can be calculated as the sum of:
  - vmstat.pgscan\_direct\_dma
  - vmstat.pgscan\_direct\_normal
  - vmstat.pgscan\_direct\_movable
    - pgscan\_d="vmstat.pgscan\_direct\_dma[0] + vmstat.pgscan\_direct\_normal[0] + vmstat.pgscan\_direct\_movable[0]"
- The available part of the cache that could be freed can be calculated as the:
  - meminfo.Cached -meminfo.Shmem
    - avail\_cache="meminfo.Cached -meminfo.Shmem"



### cpuplugd summary



- CPU Hotplug memory management will NOT release page cache memory
- The CMM module needs to be loaded before starting cpuplugd
- Understand how much memory you want to allow CMM to reclaim and the rate at which you will return memory. The last thing you want is a failing memory allocation, or adverse performance impact.
- Under heavier IO load you may want more free memory available to Linux
- The goal is for Linux to dynamically return pages of memory to z/VM when not in use, and to allow the entire system to operate more efficiently
- The amount of memory required an application to run is a function of the application program code, the workload volume, and any other software added to monitor or manage the environment.
- cpuplugd does NOT plug and unplug memory (chmem), it only uses CMM
- cpuplugd does NOT add more CPUs than what you have active at boot time







### Linux on System z Device Drivers, Features, and Commands

SC33-8411-22

http://www.ibm.com/developerworks/linux/linux390/documentation\_dev.html

Using the Linux cpuplugd Daemon to manage CPU and memory resources from z/VM Linux guests

ZSW03228-USEN-00

http://www.ibm.com/developerworks/linux/li nux390/perf/tuning\_cpuhotplug.html#cpuplu gd



### References

z/VM		IBM	
CP Co Version 6 Relea	mmands and Utiliti <sup>*3</sup>	ies Reference	
	z/VM		IBM
	Directory Main Commands Ref	itenance Facility Ference	
	Version 6 Release 3		
			SC24-6189-03

http://www.vm.ibm.com/librarys H A R E

• z/VM CP Commands and Utilities Reference SC24-6175-01

• z/VM Directory Maintenance Facility Commands Reference SC24-6188-03



## Thank you for attending!



## Please remember to fill out your session evaluation



Virtualization & Linux on

System z

Tel 262 893 8662 Email: ryoung1@us.ibm.com



