



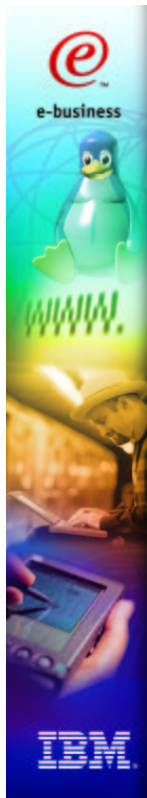
Monitoring and Understanding Performance on Linux for zSeries & S/390

SHARE 98 Technical Conference
March 3-8, 2002, Nashville
Session 5542, 9201

Oliver Benke

IBM Böblingen Lab
Schönaicher Str. 220
D-71032 Böblingen
Germany

Email: benke@de.ibm.com



Why Performance Measurement?

- Performance Tuning, Problem drill-down, Online Performance Monitoring & Analysis
- Long-term Performance Monitoring
 - ▶ Capacity Planning
 - ▶ Accounting
- Program development
 - ▶ Tracing and Profiling tools for applications or even the operating system kernel itself
- Benchmarking, Sizing
- Workload Management
 - ▶ Service Level Agreements



What can be tuned?

- CPU
- I/O
 - ▶ DASD
 - ▶ Network
 - ▶ Channels
- Memory



Forget about IDLE resources !

- A mainframe can drive most resources to their capacity limits without penalties to critical business workloads
- If one virtual server (z/OS, Linux) does not need some resources (Channel bandwidth, CPU, ...), the hardware gives it to another image ready to run
- It is like a second level of scheduling - multi-tasking in another dimension



Linux for zSeries

- Virtual Server; dynamically create and destruct Linux server using z/VM.
- Idle time of one operating system can be used by another operating system, so you are wasting less resources.
- Network - HiperSockets (iQDIO): memory speed networking to connect Linux images with other Linux images or z/OS images, leading to a client-server network in a box.



LPAR

- A mainframe can be logically partitioned
- Based on LPAR weights and on the number of logical processors, the LPAR Hypervisor allocates CPU resources to the different logical partitions
- If one LPAR has nothing to do, LPAR Hypervisor gives control to another LPAR
- z/OS IRD can influence the LPAR weights and turns logical processors online and offline



z/OS IRD

- Available with z/OS V1R2
- Linux can be part of a z/OS LPAR cluster (in contrast to OS/390)
- For Linux, only the CPU management is working
 - ▶ Adjust number of logical CPUs to reduce LPAR overhead
 - ▶ Adjust LPAR weighting factors
 - ▶ No Dynamic Channel Management (DCM) or Channel Subsystem Priority Queuing
- Does not work for IFLs



z/VM

- Second level of virtualization (or first level if machine runs in Basic mode)
- Different operating system guests can share memory, CPU and I/O resources if running under z/VM
- Especially for V=R/F guests, the performance can be fairly well, but this highly depends on the application workload



z/VM V=F,R guests

- The preferred V=R guest can use hardware facilities to execute faster. V=R guests are faster than V=F guests.
- Up to five V=F and one V=R guests (if not running z/VM under LPAR)
- All V=R,F guests must reside below the 2 GB line (z/VM 4.2)
- For each QDIO device, z/VM allocates a 8 MB shadow queue below the 2 GB line (z/VM 4.2)
- QDIO (HiperSockets) is less efficient if running under z/VM

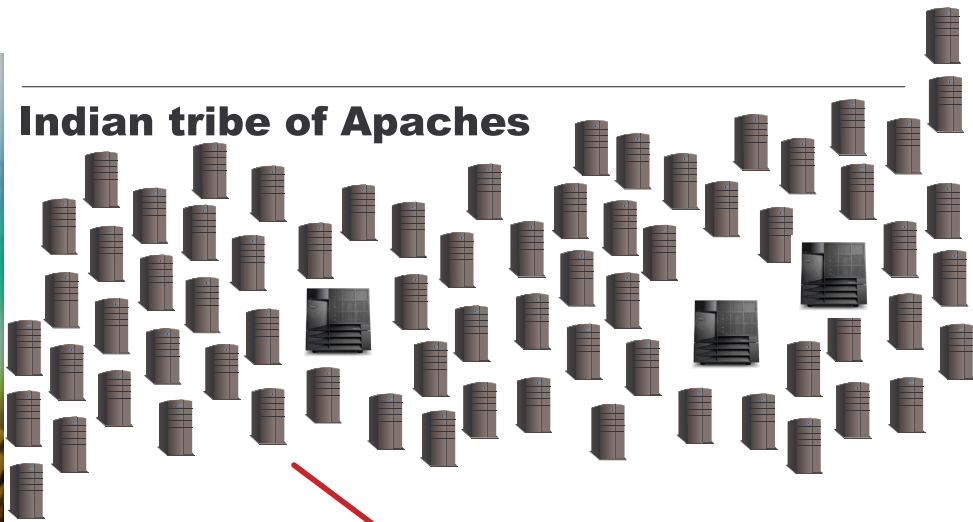


SIE

- Stands for "start interpreted execution"
- CPU instruction to facilitate virtualization layer over operating system layer as exploited by z/VM and LPAR PR/SM
- Exited for I/O processing (especially channel programs, better for HiperSockets) and if LPAR/VM time slice ends
- Allows very fast and flexible partitioning on zSeries hardware



Indian tribe of Apaches



zSeries, z/VM

"Server Farm in a Box"

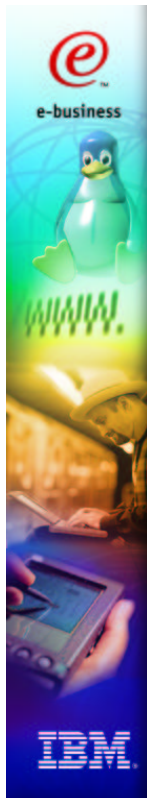
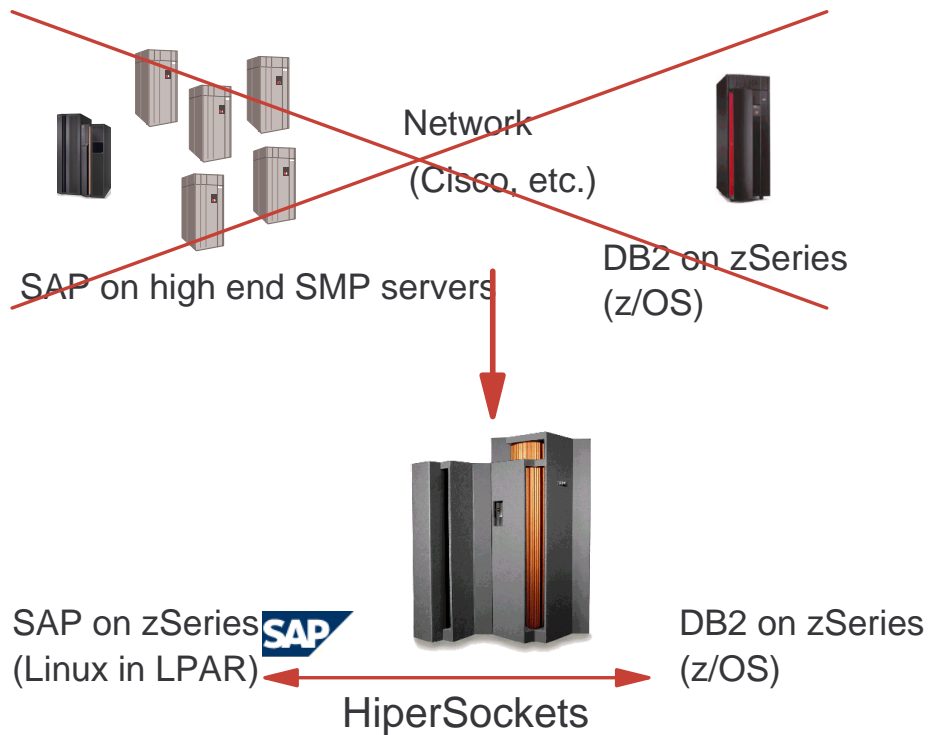


Horizontal Server Consolidation

- **Consolidate lots of under-utilized servers on one box**
 - ▶ Under-utilized web servers, mail servers, file servers, print servers
 - ▶ ISPs, ASPs or universities can give Linux servers with root access to their customers
- **For this, you definitely need z/VM**
 - ▶ Currently, LPAR is limited to 15 logical partitions per box
 - ▶ Lots of Linux images can be managed with z/VM systems management facilities



High end client/server



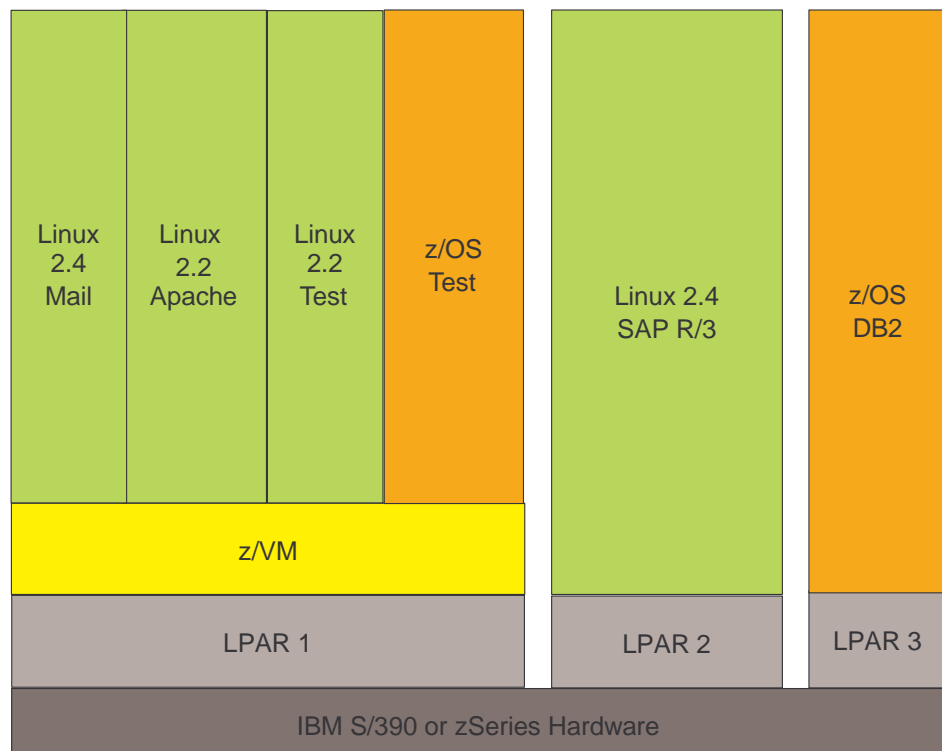
Vertical Server Consolidation

- Consolidate some high-end SMP servers on Linux for zSeries
 - ▶ WebSphere
 - ▶ SAP R/3 Application Server (together with z/OS DB2 Database Server in separate LPAR on same physical box, connected with HiperSockets)
- Probably an LPAR game
 - ▶ Faster
 - ▶ A Linux partition can be part of a z/OS LPAR cluster, so z/OS IRD can adjust LPAR weights
- Sure, you can combine horizontal and vertical server consolidation, perhaps 4 high-end virtual servers under LPAR and 1 VM LPAR for test systems and low-end server applications



Scalability of the Linux kernel

- On zSeries, Linux kernel 2.4 scales really well; you can efficiently burn all the power of a full-blown z900 with very few Linux and/or z/OS images
- Linux kernel 2.2 cannot use more than about two CPUs efficiently, even on zSeries hardware
- If you'd like to exploit Linux kernel 2.2, let z/VM do the scalability work for you: define lots of Linux operating systems scheduled and managed by z/VM





Some performance related UNIX and Linux concepts



Load average

- Average number of processes in the "run" queue
- A runnable process is one that is ready to consume CPU resources right now; a process waiting for I/O is *not* runnable
- These processes are either running or waiting for some resources to become available
- A high load average value (in relation to the number of physical processors) is an indicator for latent demand for CPU



CPU performance data reported by Linux

- You can use it for accounting if running Linux under LPAR (although LPAR CPU data obtained by a hardware interface is more precise)
- If running under z/VM, data reported by Linux can become pretty incorrect. Linux will not notice if z/VM gives all CPU resources to some other guest!



Linux Page Cache

- The page cache contains pages of memory mapped files
- It usually contains unnecessary files which can be freed, and the kernel actually discards those pages if it runs out of free memory
- On Intel Linux or for Linux running in a LPAR, the page cache is always useful as the memory would be wasted otherwise. But running under z/VM, it may cost valuable z/VM memory, leading to z/VM page activity.



Linux Buffer Cache

- A similar important Linux kernel data structure is the so-called Buffer Cache which contains pages read from or written to physical devices like DASDs
- Those pages are discarded if Linux runs out of physical memory
- Linux rarely has free space; everything not used is allocated for Page Cache and Buffer Cache, so even if Linux does not really need it all, it uses all available memory up to the last few percent.



Double Paging

- Possible for Linux under z/VM, running V=V mode (not possible for V=R,F)
- Assume page A is marked "swapped in" by Linux but paged out by z/VM; now, if Linux would like to page this page A out, first z/VM needs to page it in in order to enable Linux to page it out
- If Linux wants to page out a whole bunch of pages which were paged out previously by z/VM (not an unrealistic scenario), the system has to do a whole lot of work



Linux swap to VM virtual disc

- One solution would be to give Linux less memory and allocate a z/VM virtual disk for Linux swap space
- You can also use XPRAM (z/VM expanded storage) or a z/VM minidisk for swapping
- More details on how to efficiently use memory under z/VM are described in the ISP/ASP redbook (SG24-6299)



Linux Process memory: basic terms

- **SIZE:** size of the address space seen by the process, virtual size
- **RSS:** Resident Set Size
actual amount of memory that the process is using in RAM
- **SHARE:**
portion of the RSS that is shared with other processes, such as shared libraries



Processes and Threads

- In contrast to some commercial UNIX implementations, in Linux a thread is pretty much the same as a process, it just does not have an own address space
 - ▶ For the scheduler, a posix thread is almost like a process
 - ▶ In the /proc file system (see below), there is no difference between a process and a thread; so if you are monitoring your system, your threads might appear like processes on first sight
- As an alternative, user-space threads libraries are available

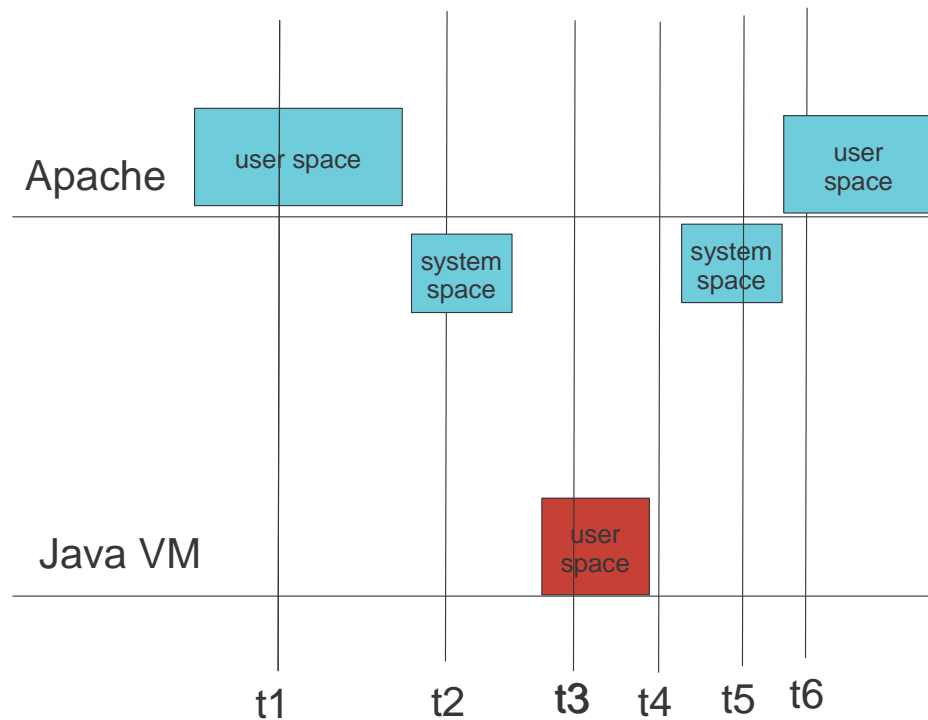


jiffies

- Derived from PC timer interrupt (100 Hz)
- Every time a timer interrupt occurs (100 times per second), the jiffies variable is incremented by one, that is one tick
- CPU usage is accounted on in units of jiffies
- If a process is running at the time the timer interrupt occurs, its CPU usage counter is incremented



When does the "jiffie event" take place?



On demand timer patch

- For an idle Linux image running under z/VM, CPU resources are used up mainly for generating the jiffies
- If you apply this patch, jiffies are generated on demand
- However, the switch between user and kernel mode is slightly slower; therefore, you should not apply this patch if running under LPAR
- see http://oss.software.ibm.com/developerworks/opensource/linux390/exp-2_4_16.shtml



Linux epoch

- The Linux scheduler works with *epochs*
- At the beginning of an epoch, each process gets a quantum for CPU resources based on its priority
- As soon as each currently runnable process is finished with its quantum, the scheduler starts a new epoch
- At the beginning of a new epoch, the dynamic priority is re-calculated based on the past behaviour of the process. The Linux scheduler tries to favour interactive processes against batch workload.



Process priorities

- Process priority can be changed with *nice/renice* commands
- Highest priority is -20, lowest priority is 19
- In addition, each process has a *dynamic priority* in Linux; a heavy CPU consumer has a worse dynamic priority than a process mainly doing I/O, giving up the CPU before the end of the time slot
- The dynamic priority is re-calculated once per epoch



System log

- Linux default: `/var/log/messages`
- Most applications are writing their error messages to `/var/log/messages`
- You should monitor the system log to find out if something went really wrong.



The `/proc` filesystem

- Virtual file system
- One of the interfaces between kernel space and user space; if the user gives a command like

```
cat /proc/stat
```

the kernel executes some function to generate the needed "virtual file"
- Parts of the `/proc` filesystem are human readable
- Most performance measurement tools for Linux are based on `/proc` file system



/proc/dasd/statistics

- Only available in Linux for zSeries, kernel version 2.4
- Used in rmfpm to calculate the following metrics:
 - ▶ dasd io average response time per request (in msec)
 - ▶ dasd io average response time per sector (in msec)
 - ▶ dasd io requests per second



/proc/dasd/statistics (continued)

```

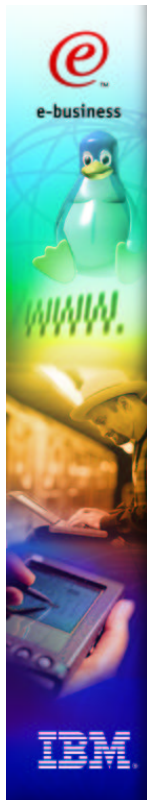
cat /proc/dasd/statistics
3156192 dasd I/O requests
__<4 __8 __16 __32 __64 __128 __256 __512 __1k __2k __4k __8k __16k __32k __64k 128k
__256 __512 __1M __2M __4M __8M __16M __32M __64M 128M 256M 512M __1G __2G __4G __>4G
Histogram of sizes (512B secs)
0 6164 0 0 0 0 0 0 0 0 0 0 0 0 0 0
Histogram of I/O times
0 0 0 0 0 0 0 0 736 628 719 952 1346 1310 448 15
4 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
Histogram of I/O times per sector
0 0 0 0 0 736 628 719 952 1346 1310 448 15 4 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
Histogram of I/O time till ssch
710 218 150 28 22 94 63 8 318 374 457 794 1271 1245 384 14
4 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
Histogram of I/O time between ssch and irq
0 0 0 0 0 0 0 0 0 3505 2072 414 147 19 2 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
Histogram of I/O time between ssch and irq per sector
0 0 0 0 0 0 3505 2072 414 147 19 2 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
Histogram of I/O time between irq and end
3 1199 959 3817 132 12 7 4 3 5 6 6 3 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0

```



/proc/stat

```
$ > cat /proc/stat
cpu 58975 2084 34136 158972653
cpu0 7792 1064 15454 26486998
cpu1 32631 993 15340 26462344
cpu2 17308 27 2320 26491653
cpu3 1240 0 614 26509454
cpu4 4 0 300 26511004
cpu5 0 0 108 26511200
page 188768 6603424
swap 0 0
intr 0
disk_io:
ctxt 1781988
btime 1011713660
processes 9867
```



/proc/slabinfo

■ statistics for frequently used kernel objects

```
cat /proc/slabinfo
slabinfo - version: 1.1 (SMP)
kmem_cache      68      68      232      4      4      1 : 252 126
nfs_read_data   0      0      384      0      0      1 : 0 62
nfs_write_data  0      0      400      0      0      1 : 0 62
nfs_page        0      0      80      0      0      1 : 0 126
tcp_tw_bucket   1      40      96      1      1      1 : 0 126
tcp_bind_bucket 136     203     16      1      1      1 : 0 126
tcp_open_request 59      59      64      1      1      1 : 0 126
inet_peer_cache 0      0      48      0      0      1 : 0 126
ip_fib_hash     8      203     16      1      1      1 : 0 126
ip_dst_cache    50      72     160      3      3      1 : 0 126
arp_cache       1      70     112      1      2      1 : 0 126
blkdev_requests 768     800     96      20     20     1 : 0 126
dnotify_cache   0      0      20      0      0      1 : 0 126
file lock cache 173     240     96      5      6      1 : 0 126
fasync cache    0      0      16      0      0      1 : 0 126
uid_cache       3      113     32      1      1      1 : 252 126
skbuff_head_cache 132     405     144     14     15     1 : 252 126
sock            85      90     816     17     18     1 : 124 62
inode_cache     28776   30296   464   3787   3787   1 : 124 62
bdev_cache      3      78      48      1      1      1 : 252 126
sigqueue       176     203     132      7      7      1 : 252 126
kiobuf          0      0     128      0      0      1 : 252 126
ccwcache-4096  0      0     4096      0      0      1 : 60 30
ccwcache-2048  4      10     2048      2      5      1 : 60 30
ccwcache-1024 118     128     1024     30     32     1 : 124 62
```



Trace facilities (Kernel patches)

- Take note on what was actually done directly in the kernel; generate trace data for some system activities
- Advantages:
 - ▶ High flexibility
 - ▶ Possibility to provide very accurate and efficient tools
- Drawbacks:
 - ▶ Has to be adopted and enabled by distributors (SuSE, RedHat); otherwise, those installing the patch are losing their service contract
- Example projects:
 - ▶ IBM dprobes
<http://www.ibm.com/developerworks/oss/linux/projects/dprobes/>
 - ▶ LTT (yes, it supports S/390)
<http://www.opersys.com/LTT/>



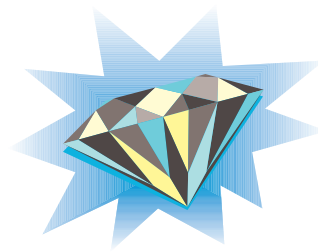
Alternative: Cycle Gatherer

- Cycle Gatherer: "Every 10 msec, make note on which processes are currently running on each of the CPUs."
- Trace Facility: "Every time the scheduler decides to switch to another process, make note on it."



Classical UNIX tools for monitoring

- sysstat package (sar, sadc)
- top
- ps
- vmstat
- free
- strace
- ...



top

- Nice option: "f - u - enter" to see what the process is waiting for

```
gfree18.boeblingen.de.ibm.com - PuTTY
8:33pm up 3:58, 3 users, load average: 1.74, 0.66, 0.24
50 processes: 47 sleeping, 3 running, 0 zombie, 0 stopped
CPU0 states: 58.0% user, 3.1% system, 0.0% nice, 38.2% idle
CPU1 states: 99.3% user, 0.2% system, 0.0% nice, 0.0% idle
CPU2 states: 1.2% user, 0.4% system, 0.0% nice, 97.3% idle
CPU3 states: 37.3% user, 0.2% system, 0.0% nice, 61.4% idle
CPU4 states: 0.0% user, 0.0% system, 0.0% nice, 100.0% idle
CPU5 states: 0.0% user, 0.0% system, 0.0% nice, 100.0% idle
Mem: 122772K av, 112752K used, 10020K free, 0K shrd, 40684K buff
Swap: 0K av, 0K used, 0K free, 23996K cached

  PID USER   PRI  NI  SIZE  RSS SHARE WCHAN  STAT %CPU %MEM  TIME  COMMAND
 1170 root    19   0   400   400   332          R    99.9  0.3   2.14  load
 1339 root    15   0 14844  14M  2236          R    48.1 12.0   0.02  cclplus
 1203 root    13   0  1560  1560  1280          R     1.3  1.2   0.01  top
    7 root     9   0     0     0     0      kupdate SW     0.5  0.0   0.01  kupdate
 1337 root     9   0   612   612   508      wait4  S     0.1  0.4   0.00  g++
    1 root     9   0   660   660   572      do_select S     0.0  0.5   0.03  init
    2 root     9   0     0     0     0      down_inte SW     0.0  0.0   0.00  kmcheck
    3 root     9   0     0     0     0      context_t SW     0.0  0.0   0.00  keventd
    4 root     9   0     0     0     0      kswapd  SW     0.0  0.0   0.00  kswapd
    5 root     9   0     0     0     0      kreclaimd SW     0.0  0.0   0.00  kreclai
    6 root     9   0     0     0     0      bdflush SW     0.0  0.0   0.00  bdflush
   63 root    -1  -20     0     0     0      end    SW<   0.0  0.0   0.00  mdrecov
```



ps - report process status

- common set of parameters:

`ps aux`

- single out a user:

`ps u --User apache`

```
bash-2.05# ps aux|more
USER      PID %CPU %MEM    VSZ   RSS TTY      STAT START   TIME COMMAND
root         1  0.0  0.1  1536   160 ?        S      Jan22   0:12 init
root         2  0.0  0.0      0     0 ?        SW     Jan22   0:00 [kmcheck]
root         3  0.0  0.0      0     0 ?        SW     Jan22   0:00 [keventd]
root         4  0.0  0.0      0     0 ?        SW     Jan22   0:22 [kswapd]
root         5  0.0  0.0      0     0 ?        SW     Jan22   0:00 [kreclaimd]
root         6  0.0  0.0      0     0 ?        SW     Jan22   0:00 [bdflush]
root         7  0.0  0.0      0     0 ?        SW     Jan22   1:05 [kupdated]
root        63  0.0  0.0      0     0 ?        SW<    Jan22   0:00 [mdrecoveryd]
root       248  0.0  0.0      0     0 ?        SW     Jan22   0:00 [keventd]
root       310  0.0  0.2  1732   292 ?        S      Jan22   0:12 syslogd -m 0
root       315  0.0  0.6  2088   768 ?        S      Jan22   0:00 klogd -2
rpc        325  0.0  0.0  1732   120 ?        S      Jan22   0:00 portmap
rpcuser   338  0.0  0.1  1844   140 ?        S      Jan22   0:00 rpc.statd
root       385  0.0  0.6  3180   800 ?        S      Jan22   0:00 /usr/sbin/sshd
root       401  0.0  0.4  2876   512 ?        S      Jan22   0:00 xinetd
```



The Process forest

- See process together with their parents or children with the `pstree` command

```
[root@lnxbnk1 /root]# pstree
init--apachegat
  |-bdflush
  |-clustergat
  |-crond
  |-dasdgat
  |-filegat
  |-find
  |-gengat
  |-gpmddsrv--gpmddsrv--2*[gpmddsrv]
  |-httpd--5*[httpd]
  |-inetd--in.telnetd--login--bash--xterm--bash--pstree
  |-keventd
  |-klogd
  |-kmcheck
```



time

- Find out how many CPU resources a command is taking

- Example:

```
$ > time make dep
...
72.52user 8.87system 2:03.72elapsed 65%CPU
(0avgtext+0avgdata 0maxresident)k
0inputs+0outputs (131158major+106391minor)
pagefaults 0swaps
$ >
```

elapsed: real time elapse
user: time this command (and its children) have spent in user space
sys: time spent in kernel space



"netstat -s" for detailed network statistics

```
$ > netstat -s

Ip:
  3608 total packets received
  0 forwarded connection openings
  0 incoming packets discarded
  3587 incoming packets delivered
  4080 requests sent out
Icmp: 493 segments received
  4 ICMP messages received
  0 input ICMP message failed.
  ICMP input histogram:
    echo requests: 4
  4 ICMP messages sent
  0 ICMP messages failed
  ICMP output histogram:
    ort received.
    echo replies: 4 rors
Tcp: 112 packets sent
  7 active connections openings
  0 passive connection openings
  0 failed connection attempts
  0 connection resets received
  3 connections established
  3493 segments received
  3964 segments send out
  10 segments retransmitted
  0 bad segments received.
  13 resets sent
Udp:
  111 packets received
  0 packets to unknown port received.
  0 packet receive errors
  112 packets sent
TcpExt:
  ArpFilter: 0
  TW: 6
  TWRecycled: 0
  TWKilled: 0
  PAWSPassive: 0
  PAWSActive: 0
  PAWSEstab: 0
  DelayedACKs: 71
  DelayedACKLocked: 0
  DelayedACKLost: 0
  ListenOverflows: 0
  ListenDrops: 0
  TCPPrequeued: 114
  TCPDirectCopyFromBacklog: 0
  TCPDirectCopyFromPrequeue: 3585
  TCPPrequeueDropped: 0
  TCPHPHits: 312
  TCPHPHitsToUser: 41
  TCPPureAcks: 1668
  TCPHPAcks: 283
  TCPRenoRecovery: 0
  TCPSackRecovery: 0
  TCPSACKReneging: 0
  TCPACKReorder: 0
  TCPSACKReorder: 0
  TCPRenoReorder: 0
  TCPTSReorder: 0
  TCPFullUndo: 0
  TCPPartialUndo: 0
  TCPSACKUndo: 0
  TCPLossUndo: 3
  TCPLoss: 0
```



free

- Give free memory; important is the second line, as buffer/cache memory is not really needed by Linux

```
[root@lnxben1 /root]# free
              total        used         free       shared    buffers     cached
Mem:          118092       116872          1220           0         4148      66124
-/+ buffers/cache:  46600       71492
Swap:           0             0             0
```



vmstat

- Gives information about memory, swap usage, I/O activity and CPU usage

```
bash-2.05# vmstat 1 10
procs          memory      swap          io          system          cpu
r  b  w  swpd  free  buff  cache  si  so  bi  bo  in  cs  us  sy  id
1  1  0    0 18608  4424 51516  0  0   0   4   0   1   0   0   4
0  1  0    0 17884  4912 51516  0  0  488   0   0  711   0   6  93
0  1  0    0 17224  5388 51516  0  0  476   0   0  512   0   9  90
0  1  0    0 16480  5800 51516  0  0  412 1196   0  447   1   7  93
0  1  0    0 14672  7016 51516  0  0 1220   0   0 1268   1  12  87
0  0  0    0 13832  7504 51516  0  0  484   0   0  571   1   3  97
0  1  0    0 12848  8080 51516  0  0  576   0   0  628   1   7  92
0  1  0    0 12228  8456 51544  0  0  376   0   0  480   2  14  84
0  1  0    0 11508  8932 51544  0  0  476 1260   0  530   0   6  94
0  1  0    0 10540  9568 51544  0  0  636   0   0  674   1   6  93
```



strace

■ Example:

```
strace -p 6148  
to trace all system calls by process with ID  
6148
```

■ Usage:

- ▶ As you can see what the process is doing, you may be able to tune it
- ▶ If you suspect a process to loop, you may check using strace; if the process consumes CPU but does not initiate any system call, it may be looping



Example: "strace ping <hostname>"

```
bash-2.05# strace ping lnxbenk1  
execve("/bin/ping", ["ping", "lnxbenk1"], [/* 23 vars */]) = 0  
uname({sys="Linux", node="gfree18", ...}) = 0  
brk(0) = 0x80017bd8  
open("/etc/ld.so.preload", O_RDONLY) = -1 ENOENT (No such file or  
directory)  
open("/etc/ld.so.cache", O_RDONLY) = 3  
fstat(3, {st_mode=S_IFREG|0644, st_size=31761, ...}) = 0  
mmap(NULL, 31761, PROT_READ, MAP_PRIVATE, 3, 0) = 0x2000001c000  
close(3) = 0  
open("/lib/libresolv.so.2", O_RDONLY) = 3  
read(3, "\177ELF\2\2\1\0\0\0\0\0\0\0\0\3\0\26\0\0\0\1\0\0\0"..., 1024) =  
1024  
fstat(3, {st_mode=S_IFREG|0755, st_size=95105, ...}) = 0  
mmap(NULL, 92712, PROT_READ|PROT_EXEC, MAP_PRIVATE, 3, 0) =  
0x20000024000  
mprotect(0x20000037000, 14888, PROT_NONE) = 0  
mmap(0x20000037000, 8192, PROT_READ|PROT_WRITE,  
MAP_PRIVATE|MAP_FIXED, 3, 0x12000) = 0x20000037000  
mmap(0x20000039000, 6696, PROT_READ|PROT_WRITE,  
MAP_PRIVATE|MAP_FIXED|MAP_ANONYMOUS, -1, 0) = 0x20000039000  
close(3) = 0
```




file system usage

■ df, du

```
bash-2.05# df
Filesystem            1k-blocks      Used Available Use%
Mounted on
/dev/dasd/6148/part1  2366164      1040288   1205680   47% /
bash-2.05# du | more
28      ./lost+found
6332    ./bin
32448   ./boot
0       ./dev/pty
0       ./dev/pts
0       ./dev/3270
0       ./dev/rd
0       ./dev/dasd/6148
0       ./dev/dasd/6149
0       ./dev/dasd
0       ./dev/discs
0       ./dev/loop
0       ./dev/md
0       ./dev
20      ./etc/X11/applnk/Utilities
```



inode utilization

- In UNIX, an inode is a structure containing meta data about files and directories.
- The number of inodes is limited, can be changed at filesystem creation time.
- If you are running out of inodes, you can not store anything more on this filesystem.
- Check with "df -i" command:

```
benke@tux390:/projects/home/benke > df -i
Filesystem      Inodes  IUsed  IFree IUse% Mounted on
/dev/dasdb1     601312  59034  542278  10% /
/dev/dasdc1     300960  63886  237074  21% /projects
```



BSD Accounting

- Writes one accounting record per terminated process or thread (as threads are something like processes in Linux...)
- Currently, SuSE decided to disable this feature for performance reasons
- Information provided:
 - ▶ user ID, group ID, process name
 - ▶ CPU resource consumption
 - ▶ average memory usage, page faults, swap activity
- An alternative to accounting Linux "from the inside" is accounting it "from the outside", with the aid of z/VM or z/OS performance tools



sysstat package

- Contains sar and sadc, long term data collector
- Normally, it collects data about overall system activity like CPU usage, swapping; no data about processes
- start with

```
$ > sadc 60 /var/log/sa/sa25 &
```

to let it generate one report every 60 seconds and write it in binary format to `/var/log/sa/sa25`
- <http://freshmeat.net/projects/sysstat/>



Mainframe-related Tools

- Some zSeries performance data is currently only available in z/VM or z/OS performance monitors
 - ▶ Coupling facility activity
 - ▶ LPAR partition data, VM CPU activity
 - ▶ Channel utilization (including OSA cards, HiperSockets)
- Tools like z/OS RMF PM and z/VM FCON can display Linux performance data together with z/OS or z/VM performance data

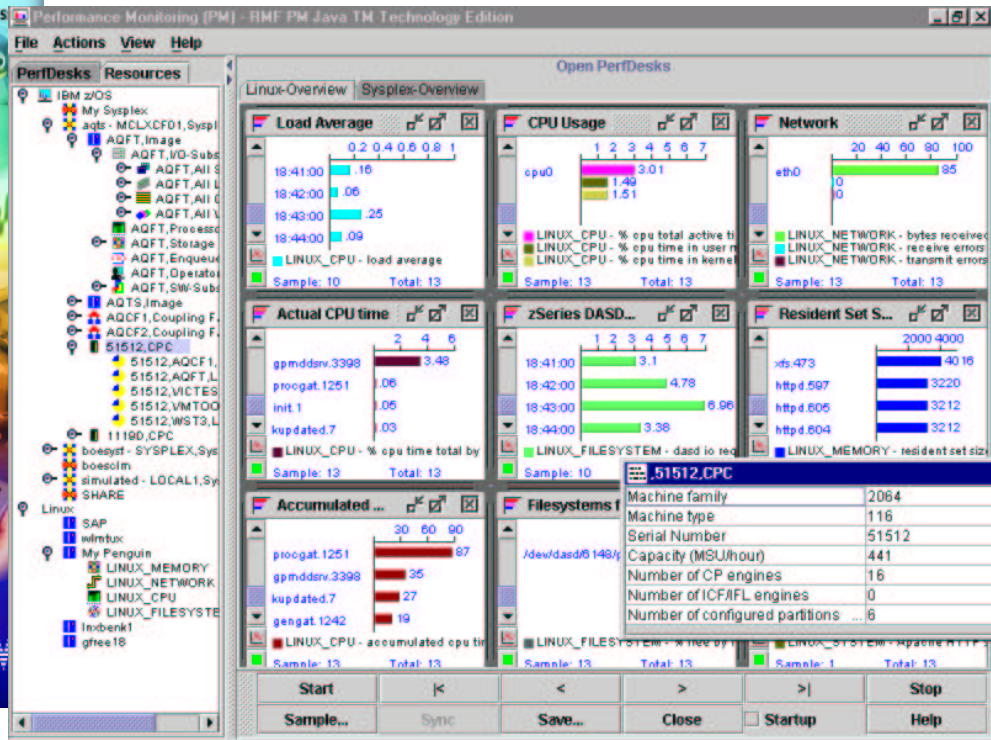


rmfpms

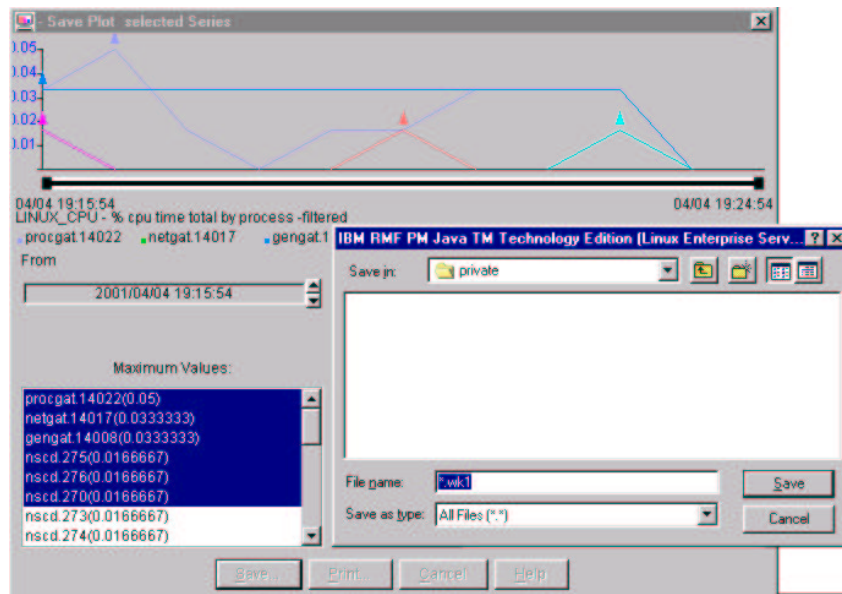
- Long term data gathering
- XML over HTTP interface
- Modular architecture
- Integrated with z/OS RMF PM and z/VM FCON
 - ▶ If you have a mixed environment with z/OS and Linux or z/VM and FCON, you can have all relevant performance metrics in one application
 - ▶ Data reported by host tools like RMF (LPAR CPU performance data, iQDIO channel utilization, etc.) is very relevant for Linux; unfortunately, we cannot make all this data available for Linux currently
 - ▶ If you have a mixed environment with z/OS, z/VM and Linux, you currently might need third-party systems management software like Tivoli DM
- see <http://www.s390.ibm.com/rmf/rmfhtmls/pmweb/pmlin.htm>



RMF PM Java Client



RMF PM: Save data in WK1 format





RMF PM Web Browser Interface

RMF DDS Browser-Interface - Microsoft Internet Explorer

Address: http://tux390.8803/

RMF DDS Browser Interface

Overview

tux390,'LINUX_CPU
% cpu time total by process

Local Time: 01/24/2002 19:08:00

in.telnetd.7401	0.05	
bash.7403	0.05	
nscd.287	0.0166667	

tux390,'LINUX_CPU
% cpu total active time by processor

Local Time: 01/24/2002 19:08:00

cpu0	4.13918	
cpu1	2.35818	

tux390,'LINUX_FILESYSTEM
% used by file system

Local Time: 01/24/2002 19:08:00

/dev/dasdc1	67.7007	
/dev/dasdb1	48.3766	

tux390,'LINUX_FILESYSTEM
size (in MB) by file system

Local Time: 01/24/2002 19:08:00

/dev/dasdc1	2310	
/dev/dasdb1	2274	

tux390,'LINUX_MEMORY
major page fault rate including children by process

Local Time: 01/24/2002 19:08:00

init.1	114	
bash.7403	64	
atd.211	0	

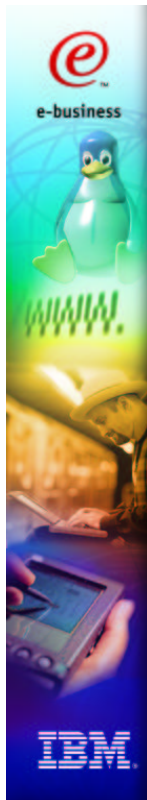
tux390,'LINUX_MEMORY
resident set size in KB by process

Local Time: 01/24/2002 19:08:00

htpd.5970	1896	
htpd.300	1828	
bash.7403	1488	

Done Local intranet

... same technology for z/OS



RMF DDS Browser-Interface - Microsoft Internet Explorer

Address: http://qgts.pok.ibm.com:8903/

RMF DDS Browser Interface

Overview

.MCLXCF01,SYSPLEX
% processor utilization by MVS image

Local Time is: 01/24/2002 13:15:00

Name	Value	Graph
AQFT	34	
AQTS	26	

.MCLXCF01,SYSPLEX
performance index by important WLM service class period

Local Time is: 01/24/2002 13:15:00

Name	Value	Graph
OMVSTASK1	3.2	
HOTTSO.2	1.3	
TSOPRIME.3	0.5	
TSOPRIME.2	0.5	

.MCLXCF01,SYSPLEX
io intensity by volume

Local Time is: 01/24/2002 13:15:00

Name	Value	Graph
AQTS.OEDEV	2354	
AQTS.C90LNH	1992	
AQTS.C90LN3	1990	
AQTS.C90BG2	1983	

.MCLXCF01,SYSPLEX
% CSA utilization by MVS image

Local Time is: 01/24/2002 13:15:00

Name	Value	Graph
AQFT	41	
AQTS	36	

Done Internet



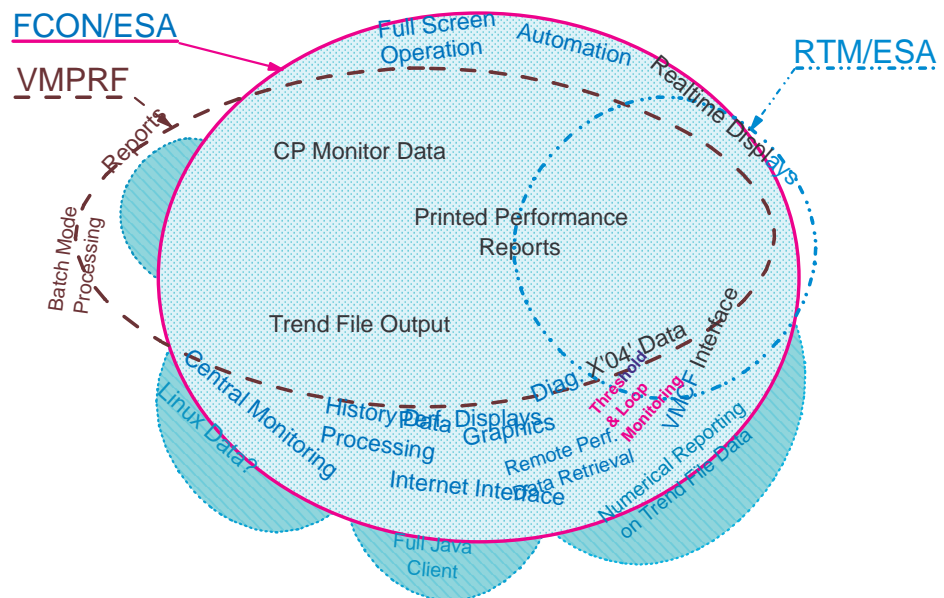
IBM FCON/ESA V.3.2.03

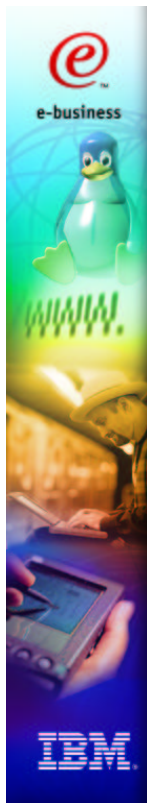
VM/ESA Full Screen Operator Console and Graphical Realtime Performance Monitor (5788-LGA) is a very powerful z/VM performance monitor. As it can display performance data collected by rmpfms in Linux, you can see VM and Linux performance data in one application.

The developer is Eginhard Jaeger (ja@ch.ibm.com), IBM Switzerland.

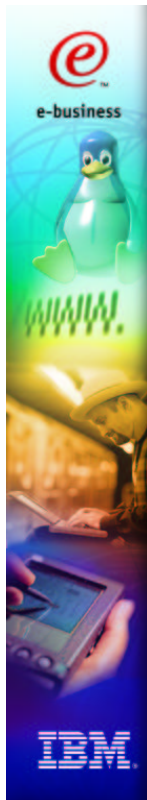
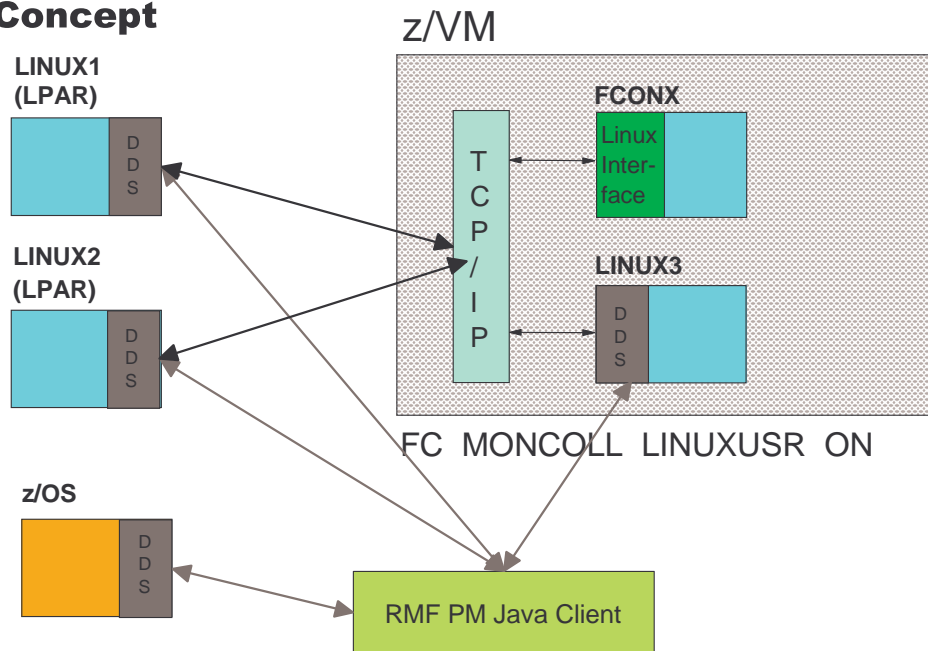


FCON: The Wishlist





Accessing Linux Performance Data Concept



Accessing Linux Perf. Data ... System Definition

File FCONX LINUXUSR

```

*****
** Initialization file with IP address definitions **
** for Linux systems that may have to be monitored. **
*****
*
LINUX1  1.111.111.111:8803
LINUX2  2.222.222.222:8803
LINUX3  3.333.333.333:8803
...
...

```

➔ Defines IP addresses of Linux systems from which performance data may have to be retrieved.
You can only monitor systems defined in this file!



LINUX Display (Linux Systems Selection Menu)

```
FCX223      CPU 9672  SER 65993      Linux Systems      Perf.
Monitor

      Selectable Linux Systems
FCONMNT      GFREE18      LNXBENKE      LNXBENK1      NC      TUX390
WLMTUX      W3PILOT1      W3VML
```

Selectable via MENU item 29 or with LINUX command



LINUX *userid* Linux Details Selection

```
FCX224      CPU 9672  SER 15585  Interval 01:28:00 - 01:29:00  Perf.
Monitor

Linux Performance Data Selection for System NC

System Data
Processes created per second      0.733
Context switches per second      52.36
Apache: Requests per second      0.017
  Bytes per request      425
  Busy threads      1
  Idle threads      4
  404 Errors per minute      0

S Perform. Reports      Description
_ LXCPU      NC      CPU utilization details
_ LXMEM      NC      Memory utilization & activity details
_ LXNETWRK      NC      Network activity (overall & by device)
_ LXFILSYS      NC      File system size and utilization
```




LXCPU *userid*

LINUX CPU Utilization Details

```

FCX230      CPU 9672  SER 15585  Interval 01:33:00 - 01:34:00  Perf.
Monitor

Linux CPU Utilization for System NC

(s)->
Processor          <--- Percent CPU Utilization --->  <-Accumulated
KernTm
>>Mean>>          Total  User  Kernel  Nice  Idle  TotTm  UserTm
---
cpu0                0.78  0.06  0.71   0   99.29  ---   ---
cpu1                0.73  0.18  0.54   0   99.26  ---   ---
cpu2                0.48   0    0.48   0   99.51  ---   ---
cpu3                0.86  0.01  0.84   0   99.13  ---   ---

Process Name
syslogd.293        0.78   0    0.78   0   ---   3657  36.08
3621
nmbd.499           0.46  0.03  0.43   0   ---   9166  1649
7517
apachegat.29502    0.3   ...   0.3   ...  ---   353.9  1.48
352.5
gengat.29511       0.3   ...   0.3   ...  ---   498.0  1.37
496.7
procgat.29517      0.23  ...   0.23  ...  ---   370.7  13.47
357.2
httpd.464          0.18  0    0.18   0   ---   953.6  18.3
935.2
  
```



LXMEM *userid*

LINUX Memory Utilization and Activity Details

```

FCX229      CPU 9672  SER 15585  Interval 01:38:00 - 01:39:00  Perf.
Monitor

Linux Memory Util. & Activity Details for System NC

Total memory size      1961MB  Swap space size      2047MB
Total memory used      1708MB  % Swap space used    0%
Used for buffer         1528MB  Swap-in rate         0/s
Used for shared         30MB     Swap-out rate        0/s
Used for cache          12MB     Page-in rate         4.783/s
Total free memory      252MB     Page-out rate        4.783/s

----->
<----- Size ----->  <----- Page Fault Rate/s
(Bytes)      (kB)      Minor  Major
<-Incl.Children->
Process Name  VirtSize  ResidSet  MinPgFlt  MajPgFlt  MinPFfltC
MajPFfltC
gpmddsrv.1904  4907010  2404     ....     ....     1
20
gpmddsrv.1908  4907010  2404     ....     ....     ....
.....
httpd.11108    2641920  1408     0         0         0
httpd.14101    2641920  1408     0         0         0
httpd.1681     2641920  1408     0         0         0
httpd.19085    2641920  1408     0         0         0
...
  
```



LXNETWRK *userid*

LINUX Network Activity

FCX227 CPU 9672 SER 15585 Interval 01:41:00 - 01:42:00 Perf.
Monitor

Linux Network Activity for System NC

Network Device	Received/s RcvPack RcvByte RcvError	Transmitted/s SndPack SndByte SndError
>Total>	45.01 5238 0	2.11 732 0
lo	0.18 13 0	0.18 13 0
tr0	44.83 5224 0	1.93 718 0



LXFILSYS *userid*

LINUX Filesystem Usage

FCX228 CPU 9672 SER 15585 Interval 01:43:00 - 01:44:00 Perf.
Monitor

Linux Filesystem Usage for System NC

DASD I/O Activity

I/O request rate per second 3.98
I/O response time/request (msec) 0.47
I/O response time/sector (msec) 0.059

Filesystem Name	Size	Free	%Used	%Free
>Total>	48204	3890	91.4	8.5
/dev/sda1	249	33	85.9	14.0
/dev/sda10	998	99	89.5	10.4
/dev/sda11	9449	1241	86.1	13.8
/dev/sda5	9868	682	92.7	7.2
/dev/sda6	9868	566	93.9	6.0
/dev/sda7	9868	347	96.2	3.7
/dev/sda8	7904	922	87.6	12.3



VM FCON SYSSUM Display

System Summary Log

```

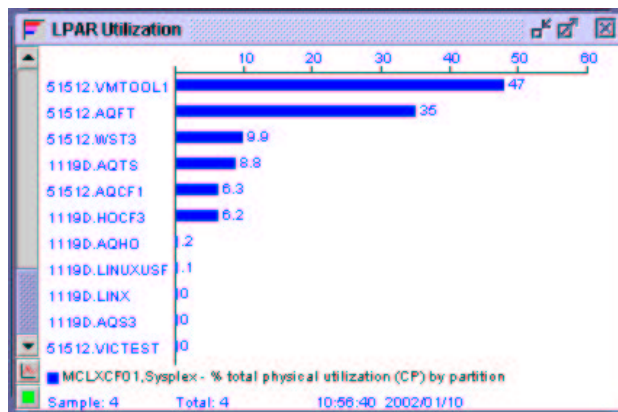
FCX225      CPU 2064  SER 51524  Interval 00:00:24 - 05:13:24  Perf.
Monitor

<----- CPU -----> <Vec> <--Users--> <---I/O---> <Stg>
<-Paging-
      <--Ratio-->
      SSCH  DASD  Users
<-Rate/s-
Interval  Pct      Cap-  On-  Pct  Log-      +RSCH  Resp  in  PGIN+
Rea      Busy  T/V  ture line  Busy  ged Activ  /s  msec  Elist  PGOUT
Wri
>>Mean>> .2 3.33 .5727 7.0 .0 41 19 58.7 .9 .0 .4
04:52:24 .2 3.93 .5612 7.0 .0 41 19 56.3 .6 .0 .0
04:53:24 .2 3.64 .5752 7.0 .0 41 19 58.6 .9 .0 .4
04:54:24 .2 3.82 .5607 7.0 .0 41 20 57.4 .6 .0 .1
04:55:24 .1 3.78 .5730 7.0 .0 41 18 56.4 .9 .0 .0
04:56:24 .2 3.09 .6105 7.0 .0 41 21 59.7 .9 .0 .0
04:57:24 .1 3.44 .6055 7.0 .0 41 20 57.1 .6 .0 .0
04:58:24 .2 3.64 .5730 7.0 .0 41 18 59.0 .9 .0 .0
04:59:24 .2 3.87 .5538 7.0 .0 41 20 57.6 .6 .0 .0
05:00:24 .2 3.74 .5699 7.0 .0 41 19 58.2 .9 .0 .0
05:01:24 .3 2.02 .7186 7.0 .0 41 20 62.0 1.0 .0 .0
05:02:24 .1 3.71 .5726 7.0 .0 41 20 58.2 .6 .0 .2
05:03:24 .2 3.87 .5611 7.0 .0 41 20 57.2 .9 .0 .0
05:04:24 .1 3.69 .5737 7.0 .0 41 20 57.1 .6 .0 .0
05:05:24 .1 3.77 .5639 7.0 .0 41 19 57.5 .9 .0 .0

```



LPAR Partition Data (from z/OS RMF)





HiperSockets display in VM FCON

FCX231 CPU 2064 SER 51524 Interval 06:55:22 - 06:56:22 Perf. Monitor

Channel Path		Hipersocket Activity/Sec.							
		<--- Total for System --->				<----- Own Partition ----->			
		<-Transferred-->		Failed		<-Transferred-->		<--- Failed --->	
ID	Shrd	T_Msgs	T_DUnits	T_NoBuff	L_Msgs	L_DUnits	L_NoBuff	L_Other	
FB	No	0	0	0	0	0	0	0	0
FC	No	0	0	0	0	0	0	0	0
FD	No	0	0	0	0	0	0	0	0
FE	No	0	0	0	0	0	0	0	0



HiperSockets Display in z/OS RMF

z/OS VIR2 SYSTEM ID CB88 DATE 07/22/2001 INTERVAL 22.54.336
 RPT VERSION VIR2 RMF TIME 15.37.05 CYCLE 1.000 SECONDS

IODF = 01 CR-DATE: 05/10/2000 CR-TIME: 21.00.01 ACT: POR MODE: LPAR CPMF: EXTENDED MODE

CHANNEL PATH ACTIVITY

OVERVIEW FOR DCM-MANAGED CHANNELS

CHANNEL GROUP	G NO	UTILIZATION(%) PART TOTAL	READ(MB/SEC) BUS PART TOTAL	WRITE(MB/SEC) PART TOTAL
FC_SM	1 8	15.36 55.86	6.00 15.36 60.00	15.36 60.36
FCV_M	12	30.00 45.00	5.00 45.00 50.00	45.00 50.00
CNC_M	1	17.23 34.45		

DETAILS FOR ALL CHANNELS

CHANNEL ID	PATH	UTILIZATION(%) G SHR PART TOTAL	READ(MB/SEC) BUS PART TOTAL	WRITE(MB/SEC) PART TOTAL	CHANNEL ID	PATH	UTILIZATION(%) G SHR PART TOTAL	READ(MB/SEC) BUS PART TOTAL	WRITE(MB/SEC) PART TOTAL
78	CVC_P	OFFLINE			80	CNC_S	OFFLINE		
79	CNC_S	OFFLINE			81	CNC_S	0.04	0.04	
7A	FC	1 Y 20.00 30.00	5.00 20.00 30.00	20.00 50.00	82	FC	Y 20.00 30.00	6.00 20.00 30.00	20.00
7B	FC_SM	Y 15.36 55.86	6.00 15.36 60.00	15.36 60.36	83	FC	1 Y 15.36 55.66	7.00 15.36 60.00	15.36
7C	FCV_M	Y 10.00 30.00	5.00 10.00 50.00	10.00 50.00	84	FCV_M	Y 10.00 30.00	5.00 10.00 50.00	50.00
7D	FCV_M	Y 30.00 45.00	5.00 45.00 50.00	45.00 50.00	85	FCV_M	Y 30.00 45.00	6.00 45.00 50.00	45.00
7E	CNC_M	17.23 34.45			86	CNC_S	0.00	0.00	
7F	CNC_S	OFFLINE			8C	CNC_S	0.00	0.00	

CHANNEL ID	PATH	WRITE(B/SEC) PART TOTAL	MESSAGE RATE PART TOTAL	MESSAGE SIZE PART TOTAL	SEND FAIL PART	RECEIVE FAIL PART TOTAL
AB	IQD	Y 645.12M 2500.2G	850.23K 4.2K	760.12 779.56	12	85 120



SIZE390

- If you need a sizing for Linux for zSeries & S/390, ask your IBM business partner or sales representative for it. They should have access to a tool called SIZE390 for
 - ▶ WebSphere Commerce Suite
 - ▶ Sendmail Advanced Mail Server
 - ▶ Linux Server Consolidation
 - ▶ later Samba, Apache, WebSphere Application Server, program development



Interface between Linux kernel and z/VM CP

- CP device driver, developed by Neale Ferguson; interface between Linux and z/VM
- <http://penguinvm.princeton.edu/programs/cpint.tar.gz>
- "#cp ind user" in Linux console:

```
CP IND
AVGPROC-069% 07
XSTORE-000037/SEC MIGRATE-0000/SEC
MDC READS-000001/SEC WRITES-000000/SEC HIT RATIO-094%
STORAGE-024% PAGING-0000/SEC STEAL-000%
Q0-00071 Q1-00000 Q2-00000 EXPAN-001 Q3-00000
EXPAN-001
```



Example Scenario

- The following Linux image may be completely idle:

```
$ > top 12:30pm
```

```
up 4 min, 2 users, load average: 0.02, 0.07, 0.03
```

```
24 processes: 23 sleeping, 1 running, 0 zombie, 0 stopped
```

```
CPU0 states: 0.1% user, 19.1% system, 0.0% nice, 80.8% idle
```

```
CPU1 states: 0.0% user, 23.2% system, 0.0% nice, 76.8% idle
```

```
...
```

- ... as z/VM is heavily loaded and does not give Linux many resources, so even for simple tasks, Linux needs about 20% of its CPU resources just to do almost nothing:

```
$ > #CP IND
```

```
AVGPROC-099% 07
```

```
...
```



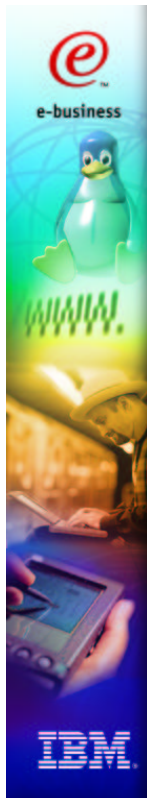
Velocity Software

- Integrates Linux and VM performance data in one application
- Uses UCD SNMP for Linux performance data:
<http://net-snmp.sourceforge.net/>
- <http://www.velocitysoftware.com>
- SHARE Session 5545/9222 (yesterday afaik)



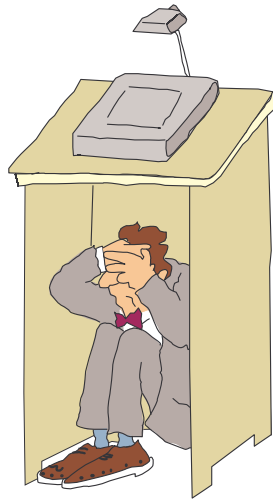
Other network monitoring tools available for Linux for zSeries & S/390

- **BigBrother** <http://www.bb4.com>
- **OpenNMS** <http://www.opennms.org>
- **MRTG**
<http://people.ee.ethz.ch/~oetiker/webtools/mrtg/>
(or <http://www.mrtg.org>)
- **NetSaint** <http://www.netsaint.org>
- **Scotty**
<http://wwwhome.cs.utwente.nl/~schoenw/scotty/>
- **BigSister** <http://bigsisiter.graeff.com/>
- ... and many more



Further Reading

- **Linux for IBM eServer zSeries and S/390:**
"*Distributions*" Redbook, SG24-6264
- **Linux for IBM eServer zSeries and S/390:**
"*ISP/ASP Solutions*" Redbook, SG24-6299
- **Jason R Fink & Matthew D Sherer:** "*Linux Performance Tuning and Capacity Planning*", SAMS 2001, ISBN 0-672-32081-9



Questions



Email: benke@de.ibm.com