# HP Escalation Team Performance Troubleshooting Techniques and Tools

**Ken Johnson**

Escalation Engineer
Hewlett Packard

# Americas Escalation Team Performance troubleshooting

Ken Johnson

Escalation Engineer

Hewlett-Packard Company

19410 Homestead Rd.  MS 4345

Cupertino, CA   95014

E-mail: ken.johnson@hp.com

Phone:  (408) 447-1307

# Americas Escalation Team Performance troubleshooting

- To share the strategies, tactics and tools used by the HP Americas Escalation Team (AET) to resolve performance escalations – using real world examples and case studies

- We will not deal with system tuning, capacity planning or benchmarking

# Outline

1. The AET Perspective: Emergency Room
2. Defining the Performance Problem
3. Metrics and Tools
4. Is there a Bottleneck ?
5. Is the Work Necessary?
6. Looking for Anomalies
7. Isolating Components
8. Knowing Your System
9. Rules of Thumb

# The AET perspective: emergency room

- ER the TV show - we do triage to stop the bleeding

- Stabilize the system as fast as possible

- Quickly identify the first steps
  - Is this an HP defect / config issue / 3$^{rd}$ party issue ?
  - Often our value-add is to point in the right direction

- We have a system perspective
  - Understand interactions between HW, OS, Network, DB, Application
  - We train our engineers for a system perspective

# Defining the performance problem

Things we want to know in the first minutes

- When did the performance  problem start ?
- How do you know you have a problem ?
  - Is this a user/business impacting problem or a metric-only issue ?
  - Is the problem quantified ?
- Is the hardware and OS base stable and consistent ?
- What changed ?

# Defining the performance problem

- Quantification
  - Allows you to measure the objective effect of changes
  - Define the current state and the goal

- Changing only one thing at a time

- Characterize and focus the problem
  - System wide or particular application?
  - All the time or specific time of day?
  - Network access or local access?
  - NFS mounts or local disks?
  - Consistent or erratic?

# Metrics & Tools

- What are the thermometers really measuring

- Metrics are simply statistics produced by software
  - Some of our escalations are with performance tools

- Always have more than one data point and always use more than one tool

- Your tools can affect the environment

# Metrics & Tools

Be sure what a metric is really measuring

- Wait time/service time

- Page out/swap out/deactivation

- Run queue/load average

- Inode table utilization

# Metrics & Tools
## first tools

- uptime

- model

- uname –a

- sysdef

- top

# Metrics & Tools
## first metric: system/user cpu

- What is system CPU?

- Why is it important?
  - Points at initial directions to pursue root cause
  - HP owns this code

- High system CPU can point to:
  - High number of system calls
  - Memory I/O problems
  - Thrashing and spinning in the kernel

- CMT has visibility into system CPU utilization
  - There are utilities we use to do kernel profiling on production systems

# Metrics & Tools
## sar – cpu report

```
$ sar 5 5

HP-UX karoo B.11.00 A 9000/820     05/09/03

17:47:15     %usr      %sys      %wio      %idle
17:47:20       12        23         7         58
17:47:25        1         6         6         87
17:47:30        0         0         1         98
17:47:35        0         0         2         98
17:47:40        0         4         1         94

Average         3         7         3         87
```

# Metrics & Tools
## sar – disk report

```
# sar -d 1 2


HP-UX cecl3 B.11.00 U 9000/800     07/08/03


14:33:20   device    %busy  avque  r+w/s  blks/s  avwait  avserv
14:33:21   c1t2d0    79.21  1.48    208    2259    6.05    7.33
14:33:22   c1t2d0    77.00  0.93    272    2244    5.53    4.21
           c1t0d0     2.00  0.50      2       4    1.37    6.17


Average    c1t2d0    78.11  1.17    240    2252    5.76    5.57
Average    c1t0d0     1.00  0.50      1       2    1.37    6.17
```

# Metrics & Tools
## measureware tools

- Tightly integrated with HPUX kernel

- On most mission critical systems

- Trial version available for escalations

- Glance / GPM

- Scope / PerfView

# Metrics and Tools
## Glance

```
h - Online Help          q - exit (or e)           A - Application List
g - Process List         d - Disk Report           P - PRM Group List
a - CPU By Processor     i - IO By File System     Y - Global System Calls
c - CPU Report           u - IO By Disk            F - Process Open Files
m - Memory Report        v - IO By Logical Volume   M - Process Memory Regions
t - System Tables        N - NFS Global Activity    R - Process Resources
w - Swap Space           n - NFS By System         W - Process Wait States
B - Global Waits         l - Network By Interface   L - Process System Calls
Z - Global Threads       T - Trans Tracker         y - Renice Process
G - Process Threads      H - Alarm History         s - Select Process
I - Thread Resource      J - Thread Wait
S - Select Disk/NFS/Appl/Trans/Thread
```

```
----------------------------------------------------------------------------------------
CPU  Util   SAU                             U          |  51%    7%   56%
Disk Util   F                                          |   2%    5%   53%
Mem  Util   S              SU        UB                B|  95%   94%   95%
Swap Util   U                   UR        R            |  55%   54%   55%
----------------------------------------------------------------------------------------
```

|                       |       |       |     |              | PROCESS LIST |          |            | Users=    3 |
|-----------------------|-------|-------|-----|--------------|--------------|----------|------------|-------------|
|                       |       |       |     | User         | CPU Util     | Cum      | Disk       | Thd         |
| Process Name          | PID   | PPID  | Pri | Name ( | 200% max) | CPU      | IO Rate    | RSS     Cnt |

| Process Name | PID   | PPID  | Pri | Name | ( 200% max) | Cum CPU | Disk IO Rate | RSS   | Thd Cnt |
|--------------|-------|-------|-----|------|-------------|---------|--------------|-------|---------|
| memm3        | 24951 | 24926 | 239 | kenj | 99.8/77.6   | 12.8    | 0.0/ 0.3     | 692kb | 1       |
| glance       | 24919 | 24905 | 154 | kenj | 0.8/ 1.1    | 1.7     | 0.0/ 0.0     | 4.9mb | 1       |
| glance       | 24632 | 23393 | 154 | root | 0.8/ 0.8    | 8.1     | 0.0/ 0.0     | 5.0mb | 1       |
| midaemon     | 1303  | 1     | -16 | root | 0.2/ 0.1    | 4373.4  | 0.0/ 0.0     | 9.3mb | 3       |
| diaglogd     | 1574  | 1161  | 168 | root | 0.0/ 0.0    | 88.5    | 0.0/ 0.0     | 440kb | 1       |
| nfsd         | 1529  | 1522  | 154 | root | 0.0/ 0.0    | 0.0     | 0.0/ 0.0     | 472kb | 1       |
| nfsd         | 1530  | 1522  | 154 | root | 0.0/ 0.0    | 0.0     | 0.0/ 0.0     | 472kb | 1       |
| nfsd         | 1531  | 1522  | 154 | root | 0.0/ 0.0    | 0.0     | 0.0/ 0.0     | 472kb | 1       |
| registrar    | 1581  | 636   | 154 | root | 0.0/ 0.0    | 0.7     | 0.0/ 0.0     | 388kb | 1       |
| nfsd         | 1532  | 1522  | 154 | root | 0.0/ 0.0    | 0.0     | 0.0/ 0.0     | 472kb | 1       |
| nfsd         | 1533  | 1521  | 154 | root | 0.0/ 0.0    | 0.0     | 0.0/ 0.0     | 472kb | 1       |
| nfsd         | 1534  | 1521  | 154 | root | 0.0/ 0.0    | 0.0     | 0.0/ 0.0     | 468kb | 1       |

```
----------------------------------------------------------------------------------------
Process     CPU      Memory    Disk      37    1     Next      Select      Help      Exit
List        Report   Report    Report                Keys      Process               Glance
```

```
CPU  Util   SAU                         U          |  52%   16%   56%
Disk Util   F                                      |   2%    5%   53%
Mem  Util   S            SU        UB            B  |  95%   94%   95%
Swap Util   U                 UR        R          |  55%   54%   55%
```

## SYSTEM TABLES REPORT                                    Users=     3

| System Table | Available | Used | Utilization | High(%) |
|---|---|---|---|---|
| Proc Table (nproc) | 5620 | 143 | 3 | 3 |
| File Table (nfile) | 10539 | 640 | 6 | 6 |
| Shared Mem Table (shmmni) | 200 | 13 | 7 | 7 |
| Message Table (msgmni) | 50 | 2 | 4 | 4 |
| Semaphore Table (semmni) | 64 | 31 | 48 | 48 |
| File Locks (nflocks) | 200 | 37 | 19 | 19 |
| Pseudo Terminals (npty) | 60 | 0 | 0 | 0 |
| Buffer Headers (nbuf) | na | 72930 | na | na |

```
------------------------------------------------------------------------------------
CPU  Util   SSAU                               U                | 54%   21%   56%
Disk Util   F                                                   |  2%    4%   53%
Mem  Util   S                 SU         UB                  B  | 95%   94%   95%
Swap Util   U                       UR         R               | 55%   54%   55%
------------------------------------------------------------------------------------
```

                          SYSTEM TABLES REPORT                         Users=    3


| System Table | Available | Requested | Used | High |
|---|---|---|---|---|
| Inode Cache (ninode) | 6488 | na | 0 | 0 |
| Shared Memory | 12.5gb | 30.2mb | | |
| Message Buffers | 800kb | na | 0kb | 0kb |
| Buffer Cache | 512.0mb | na | 512.0mb | na |
| Buffer Cache Min | 51.2mb | | | |
| Buffer Cache Max | 512.0mb | | | |
| DNLC Cache | 11608 | | | |


```
Model    : 9000/820/D380    Phys Memory : 1024mb  Network Interfaces :     4
OS Name  : HP-UX            Number CPUs :    2    Number Swap Areas  :     2
OS Release: B.11.11         Number Disks:   11    Avail Volume Groups:     4
OS Kernel Type: 64 bits     Mem Region Max Page Size: 64.0mb
```

```
------------------------------------------------------------------------------------
Process    CPU      Memory    Disk     37    1    Next     Select     Help      Exit
 List     Report    Report   Report              Keys     Process             Glance
```

```
---------------------------------------------------------------------------------------
CPU  Util  |SAU                          U            |   51%   30%   56%
Disk Util  |                                          |    0%    4%   53%
Mem  Util  |S              SU        UB                      B|   95%   95%   95%
Swap Util  |U                    UR        R          |   55%   54%   55%
---------------------------------------------------------------------------------------
```

                          GLOBAL SYSTEM CALLS                        Users=      3

| System Call Name | ID | Count | Rate | CPU Time | Cum CPU |
|---|---|---|---|---|---|
| read | 3 | 12 | 2.3 | 0.00066 | 0.00145 |
| write | 4 | 161 | 30.9 | 0.00441 | 0.01187 |
| open | 5 | 4 | 0.7 | 0.00061 | 0.00168 |
| close | 6 | 4 | 0.7 | 0.00061 | 0.00162 |
| time | 13 | 634 | 121.9 | 0.00132 | 0.00333 |
| brk | 17 | 0 | 0.0 | 0.00000 | 0.00004 |
| lseek | 19 | 7 | 1.3 | 0.00003 | 0.00006 |
| getuid | 24 | 0 | 0.0 | 0.00000 | 0.00000 |
| alarm | 27 | 0 | 0.0 | 0.00000 | 0.00002 |
| access | 33 | 1 | 0.1 | 0.00005 | 0.00025 |
| stat | 38 | 17 | 3.2 | 0.00124 | 0.00982 |

          Cumulative Interval:        11 secs

| Global | Global | System | 37 | 1 | Next | Netwk By | NFS | NFS By |
| Waits | Syscalls | Tables | | | Keys | Intrface | Global | System |

```
CPU  Util  SAU                              U              |  51%   32%   56%
Disk Util  F                                               |   1%    3%   53%
Mem  Util  S             SU        UB                   B  |  95%   95%   95%
Swap Util  U                  UR         R                 |  55%   54%   55%
```

                              GLOBAL WAIT STATES                       Users=    3

|            |      |        | Procs/  |            |      |        | Procs/  |
|            | %    | Time   | Threads | Blocked On | %    | Time   | Threads |
|------------|------|--------|---------|------------|------|--------|---------|
| IPC        | 0.7  | 10.48  | 2.0     | Cache      | 0.0  | 0.00   | 0.0     |
| Job Control| 0.0  | 0.00   | 0.0     | CDROM IO   | 0.0  | 0.00   | 0.0     |
| Message    | 0.0  | 0.00   | 0.0     | Disk IO    | 0.0  | 0.00   | 0.0     |
| Pipe       | 0.3  | 5.24   | 1.0     | Graphics   | 0.0  | 0.00   | 0.0     |
| RPC        | 0.0  | 0.00   | 0.0     | Inode      | 0.0  | 0.00   | 0.0     |
| Semaphore  | 0.3  | 5.25   | 1.0     | IO         | 0.6  | 9.80   | 1.9     |
| Sleep      | 32.3 | 498.18 | 94.9    | LAN        | 0.0  | 0.00   | 0.0     |
| Socket     | 2.0  | 31.45  | 6.0     | NFS        | 0.0  | 0.00   | 0.0     |
| Stream     | 3.0  | 47.09  | 9.0     | Priority   | 0.0  | 0.07   | 0.0     |
| Terminal   | 0.3  | 5.24   | 1.0     | System     | 42.8 | 660.81 | 125.9   |
| Other      | 17.2 | 265.57 | 50.6    | Virtual Mem| 0.0  | 0.00   | 0.0     |

```
Global    Global    System      37    1    Next    Netwk By   NFS      NFS By
Waits     Syscalls  Tables                 Keys    Intrface   Global   System
```

# Metrics and Tools
## Scope/Perfview

- Scopeux logs data from the MI database once a minute

- The data is summarized at 5 minute intervals

- Good for trend analysis

- Can help with focused troubleshooting

☐ Enable Filters    ☐ Enable Highlights

| Process Name | Date and Time | System CPU % | User CPU % | CPU % | Phys IO Rt | IO Byte Rate | Stop Reason |
|---|---|---|---|---|---|---|---|
| PSAPPSRV | Mon Sep 2 18:45:00 2002 | 0.06 | 8.00 | 8.05 | 3.7 | 72.0 | MESG |
| PSRUN | Mon Sep 2 18:45:00 2002 | 1.07 | 8.00 | 9.08 | 0.7 | 32.1 | |
| autocons | Mon Sep 2 18:45:00 2002 | 16.73 | 6.00 | 22.73 | 0.0 | 0.0 | SOCKT |
| midaemon | Mon Sep 2 18:45:00 2002 | 4.19 | 0.00 | 12.25 | 0.0 | 0.0 | SYSTM |
| ora_arc0_ASRPRD8 | Mon Sep 2 18:45:00 2002 | 0.01 | 0.00 | 0.01 | 11.4 | 394.1 | IO |
| ora_ckpt_ASRPRD8 | Mon Sep 2 18:45:00 2002 | 1.11 | 79.00 | 80.11 | 0.9 | 7.6 | PRI |
| ora_dbw0_ASRPRD8 | Mon Sep 2 18:45:00 2002 | 0.05 | 0.00 | 0.05 | 29.3 | 234.9 | SEM |
| ora_dbw1_ASRPRD8 | Mon Sep 2 18:45:00 2002 | 0.06 | 0.00 | 0.06 | 32.1 | 258.4 | SEM |
| ora_dbw2_ASRPRD8 | Mon Sep 2 18:45:00 2002 | 0.03 | 0.00 | 0.03 | 24.6 | 197.5 | SEM |
| ora_dbw3_ASRPRD8 | Mon Sep 2 18:45:00 2002 | 0.05 | 0.00 | 0.05 | 22.8 | 186.2 | SEM |
| ora_lgwr_ASRPRD8 | Mon Sep 2 18:45:00 2002 | 1.11 | 0.00 | 1.11 | 31.3 | 679.0 | IO |
| ora_lgwr_ASRSCHD | Mon Sep 2 18:45:00 2002 | 0.01 | 0.00 | 0.01 | 5.6 | 49.3 | SEM |
| oracleASRPRD8 | Mon Sep 2 18:45:00 2002 | 0.21 | 6.00 | 6.21 | 2.0 | 45.8 | SEM |
| oracleASRPRD8 | Mon Sep 2 18:45:00 2002 | 0.09 | 7.00 | 7.09 | 7.8 | 212.0 | SEM |
| oracleASRPRD8 | Mon Sep 2 18:45:00 2002 | 7.28 | 7.00 | 14.28 | 218.4 | 5120.0 | SEM |
| oracleASRPRD8 | Mon Sep 2 18:45:00 2002 | 0.08 | 6.00 | 6.08 | 0.4 | 3.8 | SEM |
| oracleASRPRD8 | Mon Sep 2 18:45:00 2002 | 0.09 | 7.00 | 7.09 | 5.5 | 162.8 | SEM |
| oracleASRPRD8 | Mon Sep 2 18:45:00 2002 | 0.14 | 6.00 | 6.14 | 0.1 | 1.2 | SEM |
| oracleASRPRD8 | Mon Sep 2 18:45:00 2002 | 0.12 | 12.00 | 12.12 | 15.3 | 488.8 | SEM |
| oracleASRPRD8 | Mon Sep 2 18:45:00 2002 | 0.04 | 5.00 | 5.04 | 6.1 | 222.8 | SEM |
| oracleASRPRD8 | Mon Sep 2 18:45:00 2002 | 0.08 | 7.00 | 7.08 | 1.6 | 48.0 | CACHE |
| oracleASRPRD8 | Mon Sep 2 18:45:00 2002 | 0.00 | 0.00 | 0.00 | 8.2 | 80.6 | |
| oracleASRPRD8 | Mon Sep 2 18:45:00 2002 | 0.08 | 7.00 | 7.08 | 1.9 | 59.3 | SEM |
| oracleASRPRD8 | Mon Sep 2 18:45:00 2002 | 0.13 | 12.00 | 12.13 | 8.1 | 275.1 | SEM |
| oracleASRPRD8 | Mon Sep 2 18:45:00 2002 | 1.16 | 10.00 | 11.17 | 17.9 | 682.2 | SEM |
| oracleASRPRD8 | Mon Sep 2 18:45:00 2002 | 0.13 | 7.00 | 7.13 | 15.7 | 518.1 | SEM |
| oracleASRPRD8 | Mon Sep 2 18:45:00 2002 | 0.13 | 9.00 | 9.13 | 6.8 | 148.0 | PRI |
| oracleASRPRD8 | Mon Sep 2 18:45:00 2002 | 0.09 | 7.00 | 7.09 | 5.3 | 194.9 | SEM |
| oracleASRPRD8 | Mon Sep 2 18:45:00 2002 | 6.56 | 66.00 | 72.56 | 0.4 | 3.7 | SLEEP |

Select the metrics you wish to display in this window.

?

- [ ] PROC_PROC_NAME
- [ ] PROC_CPU_SYS_MODE_UTIL
- [ ] PROC_CPU_USER_MODE_UTIL
- [ ] PROC_CPU_TOTAL_UTIL
- [ ] PROC_DISK_PHYS_IO_RATE
- [ ] PROC_DISK_LOGL_READ_RATE
- [ ] PROC_DISK_LOGL_WRITE_RATE
- [ ] PROC_IO_BYTE_RATE
- [ ] PROC_STOP_REASON
- [ ] INTERVAL
- [ ] PROC_APP_ID
- [ ] PROC_CPU_CSWITCH_TIME
- [ ] PROC_CPU_CSWITCH_UTIL
- [ ] PROC_CPU_INTERRUPT_TIME
- [ ] PROC_CPU_INTERRUPT_UTIL

- [ ] PROC_DISK_FS_IO_RATE
- [ ] PROC_DISK_FS_READ
- [ ] PROC_DISK_FS_READ_RATE
- [ ] PROC_DISK_FS_WRITE
- [ ] PROC_DISK_FS_WRITE_RATE
- [ ] PROC_DISK_LOGL_IO_CUM
- [ ] PROC_DISK_LOGL_IO_RATE_CUM
- [ ] PROC_DISK_LOGL_READ
- [ ] PROC_DISK_LOGL_WRITE
- [ ] PROC_DISK_PHYS_IO
- [ ] PROC_DISK_PHYS_IO_CUM
- [ ] PROC_DISK_PHYS_IO_RATE_CUM
- [ ] PROC_DISK_SUBSYSTEM_WAIT_PCT
- [ ] PROC_DISK_SUBSYSTEM_WAIT_TIME
- [ ] PROC_DISK_SYSTEM_IO

- [ ] PROC_MEM_RES
- [ ] PROC_MEM_VIRT
- [ ] PROC_MEM_WAIT_PCT
- [ ] PROC_MEM_WAIT_TIME
- [ ] PROC_MINOR_FAULT
- [ ] PROC_NFS_WAIT_PCT
- [ ] PROC_NFS_WAIT_TIME
- [ ] PROC_OTHER_IO_WAIT_PCT
- [ ] PROC_OTHER_IO_WAIT_TIME
- [ ] PROC_OTHER_WAIT_PCT
- [ ] PROC_OTHER_WAIT_TIME
- [ ] PROC_PARENT_PROC_ID
- [ ] PROC_PRI
- [ ] PROC_PRI_WAIT_PCT
- [ ] PROC_PRI_WAIT_TIME

| OK | Defaults | Cancel | Help |

☐ **Enable Filters**    ☐ Enable Highlights

| Date and Time | Device Name | Disk % | Req Queue | Phys IO Rt |
|---|---|---|---|---|
| Tue Jun  3 02:55:00 2003 | 1/10/0/0.97.29.19.0.5.0 | 0.04 | 0.00 | 0.0 |
| Tue Jun  3 02:55:00 2003 | 0/0/2/0.6.0 | 4.48 | 0.00 | 7.6 |
| Tue Jun  3 02:55:00 2003 | 0/0/2/1.6.0 | 3.98 | 0.00 | 7.0 |
| Tue Jun  3 02:55:00 2003 | 0/4/0/0.100.9.19.0.3.1 | 36.37 | 0.09 | 49.6 |
| Tue Jun  3 02:55:00 2003 | 0/4/0/0.100.9.19.0.3.2 | 2.78 | 0.06 | 2.9 |
| Tue Jun  3 02:55:00 2003 | 0/4/0/0.100.9.19.0.3.0 | 25.21 | 0.47 | 44.1 |
| Tue Jun  3 02:55:00 2003 | 1/10/0/0.100.9.19.0.3.4 | 0.14 | 0.04 | 0.3 |
| Tue Jun  3 02:55:00 2003 | 0/4/0/0.100.9.19.0.3.3 | 0.21 | 0.00 | 0.4 |
| Tue Jun  3 03:00:00 2003 | 0/4/0/0.100.9.19.0.3.2 | 3.06 | 0.00 | 2.5 |
| Tue Jun  3 03:00:00 2003 | 0/4/0/0.100.9.19.0.3.1 | 65.21 | 0.12 | 100.1 |
| Tue Jun  3 03:00:00 2003 | 0/4/0/0.100.9.19.0.3.0 | 74.06 | 0.18 | 139.9 |
| Tue Jun  3 03:00:00 2003 | 1/10/0/0.100.9.19.0.3.0 | 76.01 | 0.19 | 142.4 |
| Tue Jun  3 03:00:00 2003 | 0/4/0/0.97.29.19.0.5.0 | 0.08 | 0.00 | 0.0 |
| Tue Jun  3 03:00:00 2003 | 0/4/0/0.100.9.19.0.3.4 | 0.19 | 0.00 | 0.3 |
| Tue Jun  3 03:00:00 2003 | 1/10/0/0.100.9.19.0.3.2 | 3.35 | 0.01 | 2.7 |
| Tue Jun  3 03:00:00 2003 | 1/10/0/0.97.29.19.0.5.0 | 0.14 | 0.00 | 0.1 |
| Tue Jun  3 03:00:00 2003 | 0/0/2/1.6.0 | 4.95 | 0.06 | 8.9 |
| Tue Jun  3 03:00:00 2003 | 0/0/2/0.6.0 | 5.95 | 0.06 | 10.0 |
| Tue Jun  3 03:00:00 2003 | 1/10/0/0.100.9.19.0.3.1 | 66.57 | 0.10 | 102.0 |
| Tue Jun  3 03:00:00 2003 | 0/4/0/0.100.9.19.0.3.3 | 2.95 | 1.88 | 3.4 |
| Tue Jun  3 03:00:00 2003 | 1/10/0/0.100.9.19.0.3.4 | 0.16 | 0.00 | 0.3 |
| Tue Jun  3 03:00:00 2003 | 1/10/0/0.100.9.19.0.3.3 | 2.60 | 1.85 | 3.4 |
| Tue Jun  3 03:05:00 2003 | 0/4/0/0.100.9.19.0.3.2 | 14.61 | 0.00 | 80.3 |
| Tue Jun  3 03:05:00 2003 | 0/4/0/0.100.9.19.0.3.1 | 43.15 | 0.13 | 55.6 |
| Tue Jun  3 03:05:00 2003 | 0/4/0/0.100.9.19.0.3.0 | 100.00 | 0.20 | 206.9 |
| Tue Jun  3 03:05:00 2003 | 1/10/0/0.100.9.19.0.3.0 | 100.00 | 0.20 | 199.0 |
| Tue Jun  3 03:05:00 2003 | 0/4/0/0.97.29.19.0.5.0 | 0.04 | 0.00 | 0.0 |
| Tue Jun  3 03:05:00 2003 | 0/4/0/0.100.9.19.0.3.4 | 0.18 | 0.00 | 0.3 |
| Tue Jun  3 03:05:00 2003 | 1/10/0/0.100.9.19.0.3.2 | 14.65 | 0.01 | 81.2 |

Select the metrics you wish to display in this window.

- [ ] BYDSK_DEVNAME
- [ ] BYDSK_UTIL
- [ ] BYDSK_REQUEST_QUEUE
- [ ] BYDSK_PHYS_IO_RATE
- [ ] BYDSK_PHYS_READ_RATE
- [ ] BYDSK_PHYS_BYTE_RATE
- [ ] BYDSK_PHYS_BYTE
- [ ] BYDSK_AVG_SERVICE_TIME
- [ ] BYDSK_DIRNAME
- [ ] BYDSK_FS_READ
- [ ] BYDSK_FS_READ_RATE
- [ ] BYDSK_FS_WRITE

- [ ] BYDSK_FS_WRITE_RATE
- [ ] BYDSK_ID
- [ ] BYDSK_LOGL_READ
- [ ] BYDSK_LOGL_READ_RATE
- [ ] BYDSK_LOGL_WRITE
- [ ] BYDSK_LOGL_WRITE_RATE
- [ ] BYDSK_PHYS_IO
- [ ] BYDSK_PHYS_READ
- [ ] BYDSK_PHYS_READ_BYTE
- [ ] BYDSK_PHYS_READ_BYTE_RATE
- [ ] BYDSK_PHYS_WRITE
- [ ] BYDSK_PHYS_WRITE_BYTE

- [ ] BYDSK_PHYS_WRITE_BYTE_RATE
- [ ] BYDSK_PHYS_WRITE_RATE
- [ ] BYDSK_RAW_READ
- [ ] BYDSK_RAW_READ_RATE
- [ ] BYDSK_RAW_WRITE
- [ ] BYDSK_RAW_WRITE_RATE
- [ ] BYDSK_SYSTEM_IO
- [ ] BYDSK_SYSTEM_IO_RATE
- [ ] BYDSK_VM_IO
- [ ] BYDSK_VM_IO_RATE
- [ ] INTERVAL

| OK | Defaults | Cancel | Help |

# Metrics & Tools
## internal tools - tusc

- Trace Unix System Calls (tusc) - like truss on Solaris

- Traces all system calls made and signals received for a process

- Displays arguments in a symbolic way

**Tusc [options] command [args] | pid [pid …]**

**Useful options:**

| | |
|---|---|
| **-c** | Summary of syscall counts, errors and CPU time |
| **-ccc** | CPU time for every syscall |
| **-T %T** | Print a timestamp before every trace |
| **-rall** | Display read buffer for all reads |
| **-wall** | Display write buffer for all writes |
| **-f** | Follow fork()s |
| **-l** | Print thread id |
| **-p** | Print pid |
| **-n** | Print process name |

# Metrics & Tools
## internal tools - tusc

```
# tusc -o tusc.out 5263

open("/dev/telnet/", O_RDONLY,) .......... = 6
fcntl(6, F_SETFD, 1) ..................... = 0
brk(0x400c0000) .......................... = 0
open("myfile", O_RDONLY, 02624).......... ERR#2 ENOENT
getdents(6, 0x4009f328, 8192)............ = 48
stat("/dev/", 0x6fff27f0)................ = 0
close(6) ................................. = 0
poll(0x6fff37c0, 1, 0) ................... = 1
setsockopt(13, SOL_SOCKET, SO_KEEPALIVE, 0xc000000023b49898,
4) = 0
recv(13, "sqAW8BPQAAsqlexe ", 4096, 0) .. = 371
ioctl(6, FIONBIO, 0xc000000023a915b8) ... = 0
```

```
# tusc -o tusc.out -c find /etc -name fred
```

| Syscall | Seconds | Calls | Errors |
|---------|---------|-------|--------|
| exit    | 0.00    | 1     |        |
| read    | 0.06    | 242   |        |
| write   | 0.10    | 227   |        |
| open    | 0.00    | 9     | 2      |
| close   | 0.00    | 6     |        |
| brk     | 0.00    | 4     |        |
| lseek   | 0.00    | 7     |        |
| execve  | 0.00    | 1     |        |
| umask   | 0.00    | 2     |        |
| mmap    | 0.00    | 11    | 1      |
| fstat   | 0.00    | 4     |        |
| sysconf | 0.00    | 2     |        |
| stat64  | 0.00    | 4     | 2      |
| -----   | ------  | ----- | ------ |
| Total   | 0.17    | 534   | 6      |

# Metrics & Tools
## internal tools - tusc

```
# tusc -o tusc.out -ccc mycmd myargs


<0.000285> open("/dev/telnet/", O_RDONLY,) = 6

<0.000123> fcntl(6, F_SETFD, 1) .......... = 0

<0.000296> brk(0x400c0000) .............. = 0

<0.000199> open("myfile", O_RDONLY, 02624).. ERR#2 ENOENT

<0.000309> getdents(6, 0x4009f328, 8192).. = 48

<0.000223> stat("/dev/", 0x6fff27f0)...... = 0

<0.000266> close(6) ..................... = 0

<0.000218> poll(0x6fff37c0, 1, 0) ....... = 1

<0.000118> setsockopt(13, SOL_SOCKET, SO_KEEPALIVE,
0xc000000023b49898, 4) = 0

<0.000165> recv(13, "sqAW8BPQAAsqlexe ", 4096, 0) = 371

<0.000189> ioctl(6, FIONBIO, 0xc000000023a915b8) = 0
```

# Metrics & Tools
## internal tools - tusc

```
#tusc  -o <filename>  -ccc -f -l -n -p -v -T "%H:%M:%S" <pid>

10:17:25  connect(53, 0x7f0025f0, 16) ... ERR#245 EINPROGRESS
                         sin_family: AF_INET
                           sin_port: 3206
                  sin_addr.s_addr: 172.18.187.8


10:17:25  poll(0x415532dc, 122, 5000).................. = 1


                      poll[52].fd: 52
                  poll[52].events: POLLOUT
                 poll[52].revents: POLLOUT


                      poll[53].fd: 53
                  poll[53].events: POLLIN|POLLPRI
                 poll[53].revents: 0
```

tracing tool that uses the same trace points as measureware

```
read()
11.2303027  cpu=0 seqcnt=253201877 pid=18753 ktid=19208 utid=0
err=0 ret1=8192 ret2=1223090432 syscallbeg=     0.000473
ktcsys=     0.000003 f_data=0x59d81208 dev_t=64/0x06000a
VNODE/VNFS_SPEC/VCHR
A0=0 A1=3 A2=0 A3=0x40001340 A4=0 A5=0x2000


QUEUEDONE
11.239211  cpu=1 seqcnt=241560007 pid=18753 ktid=19208 utid=0
merged=0 dev_t=31/0x071300 retries=0 qs=     0.063364
```

kshell script that parses Kitrace raw data files

```
Frequency of trace types...

Freq      Trace_type                 Percent_of_total
202344    gettimeofday               28.0
59598     BRELSE                     8.3
49687     select                     6.9
46865     pstat                      6.5
37841     SWTCH                      5.2
37841     SETRQ                      5.2
34693     RESUME                     4.8
32530     sigprocmask                4.5
30506     kill                       4.2
17158     time                       2.4
16050     read                       2.2
15786     times                      2.2
15093     write                      2.1
```

# Metrics & Tools
## internal tools - Kparse

```
Wait symbols...

Freq        Percent            Kernel_Symbol

16531       47.8%              real_nanosleep

 7151       20.7%              read_sleep

 4994       14.4%              semop

 1467        4.2%              $PIC$3

 1336        3.9%              lvmkd_daemon

  898        2.6%              ksleep_one

  831        2.4%              select

  216        0.6%              ogetblk

  182        0.5%              pm_sigwait

  162        0.5%              sleep_spinunlock

  160        0.5%              biowait

  153        0.4%              poll
```

# Metrics & Tools
## internal tools - Kparse

```
Searching BRELSE records for hotblocks..

Freq       Block             bptype        operation

4167,    bp=0x1cabc3c00,  DATA/vxfs     wr=read

3189,    bp=0x15e843480,  INDBK/vxfs    wr=read

2967,    bp=0x1cabaeb00,  INDBK/vxfs    wr=write

 405 ,   bp=0x1567f7080,  DATA/vxfs     wr=write

 405,    bp=0x12d085500,  DATA/vxfs     wr=write

 327,    bp=0x206087580,  DATA/vxfs     wr=write

 327,    bp=0x1e6b0e100,  DATA/vxfs     wr=write

 327,    bp=0x15d59f780,  DATA/vxfs     wr=write

 222,    bp=0x1f1a67080,  DIR/vxfs      wr=read

 213,    bp=0x159896400,  DIR/vxfs      wr=read
```

# Metrics & Tools
## internal tools - Prospect

Prospect is a performance analysis tool based on based on KI tracing
and Kernel Timing Clocks

System Summary
- System wide activity
- Per-CPU counters
- Per-Process summary

How each thread of each process spends its time
- Profiling, both user and system mode for every thread
- Instruction level profiling

Application profile requires no special compilation
- No source code access needed
- Very lightweight

# Metrics & Tools
## internal tools - Prospect

- Download from:   http://www.hp.com/go/prospect

- To run, start daemon:

  ```
  # prospect -P
  ```

- Then execute command under prospect:

  ```
  # prospect -V4 -e -f prospect.txt <command [args]>
  ```

- Useful options:
- **-V2**    Trace only the specified command
- **-V3**    Trace command and descendants
- **-V4**    Trace every process
- **-Vk**    System-wide kernel profile
- **-e**     Extended instruction-level tracing
- **-f**     Specify output file

# Metrics & Tools
## internal tools - Prospect

```
System time
   select           0.22    32%
   gettimeofday     0.18    25%
   sigprocmask      0.09    13%
   write            0.05     8%
   ioctl            0.05     7%
   read             0.03     5%


User time
   File-open-fail   4.41    69%   Dev=0x40000005,Inode=4078
   XTextExtents     0.41     6%   /usr/lib/X11R5/libX11.1
   memcmp           0.17     3%   /usr/lib/libc.1
   _isspace         0.13     2%   /usr/lib/libc.1
   MatchBranchHead  0.08     1%   /usr/lib/X11R5/libXt.1
   memmove          0.06     1%   /usr/lib/libc.1
```

# Metrics & Tools
## internal tools - Kgmon

- Activates kernel trace points
- Can provide flat or full profiles

| %time | seconds | calls | name |
|---|---|---|---|
| 24.6 | 133.99 | | IN_USER_MODE |
| 15.4 | 83.96 | | prod_fullgprof_intercept |
| 10.4 | 56.43 | | _mcount |
| 8.2 | 44.89 | 3105 | idle_nonpset_loop |
| 6.2 | 33.70 | 35968151 | soo_select |
| 5.7 | 30.88 | 14384885 | hpstreams_select_int2 |
| 2.5 | 13.72 | | asm_spinlock |
| 2.4 | 13.19 | | spinlock_usav |
| 1.9 | 10.48 | 126528 | pollscan |
| 1.9 | 10.41 | 14891228 | mp_socket_lock |
| 1.3 | 6.91 | | binit |
| 1.2 | 6.27 | 14891199 | sounlock |

# Metrics and Tools
## Internal tools - p4 tools

- A set of tools that are compiled with the libp4 library

- Libp4 provides a quick way to write c programs that can access kernel structures

- Developed by GSE/WTEC organization

- Development focus is on dump analysis, we use on live systems

# Metrics and Tools
## p4 tools - kmeminfo

```
Physical memory usage summary (in page/byte/percent):

Physmem            =   262144     1.0g 100%   Physical memory
  Freemem          =    14572    56.9m   6%   Free physical memory
  Used             =   247572   967.1m  94%   Used physical memory
    System         =   208176   813.2m  79%   By kernel:
      text         =     2361     9.2m   1%    text
      data         =      418     1.6m   0%    data
      bss          =      348     1.4m   0%    bss
      Static       =    16874    65.9m   6%    for text/static data
      Dynamic      =    59617   232.9m  23%    for dynamic data
      Bufcache     =   131072   512.0m  50%    for buffer cache
      Eqmem        =       19    76.0k   0%    for equiv.mapped memory
      SCmem        =      594     2.3m   0%    for critical memory
    User           =    44381   173.4m  17%   By user processes:
      Uarea        =     2456     9.6m   1%    for thread uareas
    Disowned       =        8    32.0k   0%   Disowned pages
```

# Metrics & Tools
## p4 tools - shminfo

```
Global 32-bit shared quadrants:

=================================
        Space       Start        End   Kbytes Usage
Q4 0x0f1a0000.0xc0000000-0xc0008fff      36 OTHER
Q4 0x0f1a0000.0xc0009000-0xc0009fff       4 SHMEM id=0
Q4 0x0f1a0000.0xc000a000-0xc000bfff       8 SHMEM id=2
Q4 0x0f1a0000.0xc000c000-0xc00e0fff     852 OTHER
Q4 0x0f1a0000.0xc00e1000-0xc00f0fff      64 SHMEM id=1 locked
Q4 0x0f1a0000.0xc00f1000-0xc02c8fff    1888 OTHER
Q4 0x0f1a0000.0xc02c9000-0xc02cefff      24 SHMEM id=407
Q4 0x0f1a0000.0xc02cf000-0xc034efff     512 OTHER


Limits for 32-bit SHMEM allocation:

=====================================

Maximum shmem segment:        65536 Kbytes (shmmax)
Largest free  segment:      1048576 Kbytes (Window id 0 (global)
quadrant 2)
Available swap      :        772492 Kbytes (swapspc_cnt)
Available pseudo-swap:        20616 Kbytes (swapmem_cnt)
```

```
        semmni      64        # of semaphore identifiers
        semmns     128        # of semaphores in system
        semmnu      30        # of undo structures in system
        Semmsl    2048        max # of semaphores per ID
        semopm     500        max # of operations per semop call
        semume      10        max # of undo entries per process
        semusz     104        size in bytes of undo structure
        semvmx   32767        semaphore maximum value
        semaem   16384        adjust on exit max value

Pending semaphore operations:

kthread at 0x4209f040 sleeping in semop():
        cmd   = "ntl_reader"
        proc  = 0x4209e040 (pid 480)
        wchan = 0xbbdae4 (sem 0xbbdae0, n-waiter)
        semid = 2
        semundo = 0xc3cd48
        struct sembuf ops[1] at 0xd4e1800.0x7f7f0620:
                op sem_num sem_op sem_flg
```

- Developed to debug Service Guard cmcld hangs

- We now use it as a trigger for any short-term hang issue

- Original design assigned a process to each cpu, then it would report delays

- Buddy system assigns a process to keep track of another process and report the delay

- Large amount of data collection options

# Metrics and Tools
## internal tools – timer9

-r make process realtime (-32 (strongest) to 127 (weakest))
-t make process timeshare (not realtime)
-m memory lock process
-s how many secs between process wakeups - default 0.5
-c if process has not run for this many secs then report – default 1
-l write to logfile instead of stdout/stderr
-v filename containing list of kernel variables
-a report all processes which have used at least percent of CPU
-k toggle kernel profiling every secs seconds
-b do "buddy" -k,-p,-a and -P handling on MP systems
-C cause box to crash (panic) when delay of >= secs occurs
-B make all other CPUs check cpunum for -b and -C options
-P run program after a delay has occcured
-p send a -S's signo to pid after delay occurs
-S signo to be sent to -p's PID (default SIGUSR1)
-g ensure that there is a minimum gap of this number of secs
   between signal/exec of -p/-P process and -D delay time
   profiling (default 5).

# Is there a bottleneck ?

- This is the supply side of performance
- Easiest to look at – easiest to fix

- IO
  - Is there queuing on any drives?
  - Are there long service/wait times on any drives?
- CPU
  - Is there a significant load average?
  - Is system CPU high?
  - Are processes priority waited?
- Memory
  - Is there any paging or deactivations?
  - Is there significant swap utilization?

# Is the work necessary ?

- Is the I/O demand efficient?

- Are the CPU cycles necessary?

- Is the application efficient?

- Is the memory utilization necessary?

# Is the I/O demand efficient ?

System:          K580 4-way 11.0          Manufacturing
                 2 GB memory
                 Database server

Symptoms:        2 year installation
                 Suddenly batch jobs taking much longer to execute
                 No changes to programs or database settings
                 No system bottlenecks
                 Elevated I/O rate but no queuing and fast service times

Diagnosis:       Cost based plan had been used for key queries in DB

                 Several of the queries started doing serial I/O

# Are the CPU cycles necessary ?

System:           T600 8-way 10.20
                  2 GB memory
                  Development system
                  Compiling and source code management


Symptoms:         3 year installation
                  Recently seeing slow overall performance
                  Intermittent
                  High system CPU and high context switch rates

Diagnosis:        Files used for compiling were located in one directory

                  Large number of files and very volatile

                  Contention around the directory file itself (25 MB)

                  Spinning while waiting for shared resource caused
                            unnecessary context switching

# Is the application efficient ?

System:          N4000 4-way 11.0
                 4 GB memory
                 Web server

Symptoms:        New installation
                 Server throughput was never acceptable
                 High CPU utilization with mostly user CPU
                 Load average was reasonable and good
                         system response time

Diagnosis:       Identified large # of semop calls in bolt-on application

                 Application was in the critical path for the server

                 Allowed vendor to identify configuration problem

# Is the memory utilization necessary ?

System:             V2500 16-way 11.0
                    4 GB memory
                    Database server


Symptoms:           Memory utilization at 100%
                     High page out and deactivation rates


Diagnosis:          Default 50% buffer cache had been used
                    maxuser had been set very high – affects many other
                    kernel variables

                    Final solution was to add memory and to tune kernel
                    variables

# Looking for anomalies

- System call rates/CPU utilization

- I/O patterns
  - By device
  - By time of day
  - By process

- Wait states
  - Global and per process

# Anomalies
## system call rates / CPU util

System:              K460 4-way 10.20
                     2 GB memory
                     Legacy shell script-based application
                     Files ftp'd in, processed, then put in a directory for
                               pickup

Symptoms:            Suddenly application throughput was down
                     No changes to the application
                     System CPU way up

```
--------------------------------------------------------------------------------
CPU  Util    S                                          |    2%     2%    14%
Disk Util                                               |    0%     0%    10%
Mem  Util    S  SU                        UB  B         |   50%    49%    50%
Swap Util    UUR         R                              |   20%    20%    20%
--------------------------------------------------------------------------------
                      GLOBAL SYSTEM CALLS                      Users=      1
System Call Name         ID    Count      Rate     CPU Time   Cum CPU
--------------------------------------------------------------------------------
exit                      1        0       0.0     0.00000    0.03828
fork                      2        0       0.0     0.00000    0.02793
read                      3      392      87.1     0.00144    0.13783
write                     4      119      26.4     0.00103    0.08626
open                      5        4       0.8     0.00018    0.03305
close                     6        4       0.8     0.00012    0.00746
wait                      7        0       0.0     0.00000    0.00009
unlink                   10        0       0.0     0.00000    0.00105
chdir                    12        0       0.0     0.00000    0.00006
time                     13      199      44.2     0.00012    0.00180
brk                      17        0       0.0     0.00000    0.00162


         Cumulative Interval:        50 secs

                                                          Page 1 of 9
--------------------------------------------------------------------------------
  Global    Global     DCE      System    68    1    Next    Netwk By    NFS     NFS By
  Waits    Syscalls   Global    Tables                Keys    Intrface   Global  System
```

# Anomalies
## system call rates / CPU util

System:    K460 4-way 10.20
        2 GB memory
        Legacy shell script-based application
        Files ftp'd in, processed, then put in a directory for
           pickup

Symptoms:   Suddenly application throughput was down
        No changes to the application
        System CPU way up

Diagnosis:   vfork() was very large CPU consumer
        Identified shell script that was in a loop

# Anomalies
## IO patterns by device, time, process

System:        N4000 4-way 11.0

                  4 GB memory

                  Database server for web front-end

Symptoms:     New installation

                  System response was good

                  Unacceptable database performance

                  DB connections were short-lived

                  Analysis showed that delay was in DB disconnect

```
-----------------------------------------------------------------------------------------
CPU  Util   SUI                                    |   5%    2%   14%
Disk Util                                          |   0%    0%   22%
Mem  Util   S  SU                       UB  B      |  50%   50%   51%
Swap Util   UUR          R                         |  20%   20%   20%
-----------------------------------------------------------------------------------------

Open Files PID: 21113, netscape          PPID: 21112 euid:    101 User: kenj
                                                        Open   Open
FD  File Name                                   Type   Mode   Count      Offset
-----------------------------------------------------------------------------------------
 12 <reg,vxfs,/home,/dev/vg00/lvol4,inode:80>   reg    rd/wr    1       131072
 13 <reg,vxfs,/home,/dev/vg00/lvol4,inode:81>   reg    rd/wr    1        16384
 14 <reg,vxfs,/home,/dev/vg00/lvol4,inode:93>   reg    rd/wr    1          260
 15 <reg,vxfs,/home,/dev/vg00/lvol4,inode:83>   reg    rd/wr    1          260
 16 <fifo,pipe,inode:0>                          fifo   read     1            0
 17 <fifo,pipe,inode:0>                          fifo   write    3            0
 18 /dev/null                                    chr    write   22         1250
 19 /dev/null                                    chr    write   22         1250
 20 <reg,vxfs,/home,/dev/vg00/lvol4,inode:136>  reg    rd/wr    1          194
 21 <socket: inet,tcp,0x009f5e00>               socket rd/wr    1        16878
 22 <socket: inet,tcp,0x02387400>               socket rd/wr    1        16043
 23 <socket: inet,tcp,0x009d0800>               socket rd/wr    1        25478
                                                               Page 2 of 3
-----------------------------------------------------------------------------------------

Process    Wait      Memory    Open      68     1    Next    Process
Resource   States    Regions   Files                 Keys    Syscalls
```

## IO patterns by device, time, process

System:        N4000 4-way 11.0
               4 GB memory
               Database server for web front-end


Symptoms:      New installation
               System response was good
               Unacceptable database performance
               DB connections were short-lived
                Analysis showed delay was in DB disconnect


Diagnosis:     Used Glance to observe when user disconnected
               Found high rates of IO during disconnect
               IO was to 2 database trace files

# Anomalies
## IO patterns by device, time, process

System:            V2600 32-way
                   16 GB memory
                   Database server

Symptoms:          Application queues building up intermittantly
                   System response was good
                   sar showed average service time was ok

Diagnosis:         Used kitrace to determine there were short
                   bursts of IO causing the EMC a problem

```
IO initiation


timestamp device      ms since
                 last IO start

0.319111    40500
0.468733    40500              149
0.504507    40500               35
0.508329    40500                3
0.513627    40500                5
0.515572    40500                1
0.520509    40500                4
0.522931    40500                2
0.523979    40500                1
0.524101    40500                0
0.663621    40500              139
0.790192    40500              126
1.757730    40500              967
```

# Anomalies
## IO patterns by device, time, process

```
IO completion


timestamp    device      service time
                          (ms)

0.330263     40500       11.154
0.479900     40500       11.167
0.516062     40500       11.557    <<< burst starts at .504507
0.543143     40500       34.816
0.566754     40500       53.129
0.588817     40500       73.247
0.605655     40500       85.149
0.617827     40500       94.898
0.628228     40500      104.251
0.634265     40500      110.166
0.680130     40500       16.511
0.803631     40500       13.441
```

```
timestamp    interval      type    service_time

0.504507        25        start                  *** first burst ***
0.508329         4        start
0.513627         5        start
0.515572         2        start
0.516062         0        complete        12
0.520509         4        start
0.522931         2        start
0.523979         1        start
0.524101         0        start
0.543143        19        complete        35
0.566754        24        complete        53
0.588817        22        complete        73
0.605655        17        complete        85
0.617827        12        complete        95
0.628228        10        complete       104
0.634265         6        complete       110
```

# Anomalies
## Wait states – global and per process

System:          V2250 8-way 11.0
                 8 GB memory
                 Database server

Symptoms:        New installation
                 Slow database throughput
                 No system bottlenecks or high utilization

```
--------------------------------------------------------------------------------
CPU  Util   SUI                                            |   6%    2%   14%
Disk Util                                                  |   0%    0%   10%
Mem  Util   S  SU                              UB  B       |  50%   49%   50%
Swap Util   UUR          R                                 |  20%   20%   20%
--------------------------------------------------------------------------------
```

|                       | GLOBAL WAIT STATES |         | Users= | 1 |
|---|---|---|---|---|---|---|---|

| Event     | %    | Time   | Procs/ Threads | Blocked On  | %    | Time   | Procs/ Threads |
|-----------|------|--------|----------------|-------------|------|--------|----------------|
| IPC       | 0.0  | 0.00   | 0.0            | Cache       | 0.0  | 0.00   | 0.0            |
| Job Control | 0.0 | 0.00  | 0.0            | CDROM IO    | 0.0  | 0.00   | 0.0            |
| Message   | 0.0  | 0.00   | 0.0            | Disk IO     | 0.0  | 0.00   | 0.0            |
| Pipe      | 0.7  | 5.09   | 1.0            | Graphics    | 0.0  | 0.00   | 0.0            |
| RPC       | 0.0  | 0.00   | 0.0            | Inode       | 0.0  | 0.00   | 0.0            |
| Semaphore | 0.0  | 0.00   | 0.0            | IO          | 0.0  | 0.00   | 0.0            |
| Sleep     | 45.9 | 353.71 | 69.6           | LAN         | 0.0  | 0.00   | 0.0            |
| Socket    | 0.0  | 0.01   | 0.0            | NFS         | 0.0  | 0.00   | 0.0            |
| Stream    | 0.7  | 5.09   | 1.0            | Priority    | 0.0  | 0.09   | 0.0            |
| Terminal  | 1.3  | 10.17  | 2.0            | System      | 38.3 | 295.20 | 58.1           |
| Other     | 13.2 | 101.60 | 20.0           | Virtual Mem | 0.0  | 0.00   | 0.0            |

Page 1 of 1

# Anomalies
## Wait states – global and per process

System:         V2250 8-way 11.0
                8 GB memory
                Database server

Symptoms:       New installation
                Slow database throughput
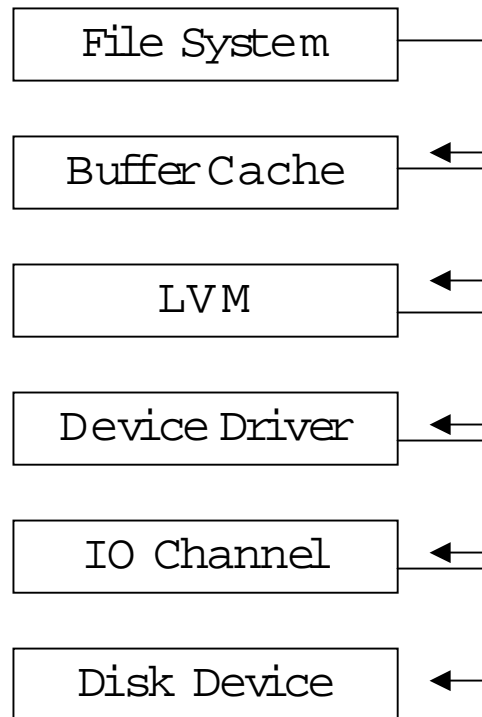                No system bottlenecks or high utilization

Diagnosis:      Identified high semop waits
                Database tuning required

# Isolating components

Make everything into a black box

Define and manipulate inputs and outputs

- The discrete elements of an I/O request

- Taking the network out of the picture

- OmniBack performance debugging techniques

# Isolating components
## discrete elements of an I/O request

```
┌─────────────────────┐
│     File System     │──────┐
└─────────────────────┘      │
┌─────────────────────┐      │
│    Buffer Cache     │◄─────┘      mincache = direct
└─────────────────────┘──────┐
┌─────────────────────┐      │
│         LVM         │◄─────┘      /dev/vg/rlvol
└─────────────────────┘──────┐
┌─────────────────────┐      │
│    Device Driver    │◄─────┘      /dev/rdsk/cxtxdx
└─────────────────────┘──────┐
┌─────────────────────┐      │
│     IO Channel      │◄─────┘
└─────────────────────┘──────┐
┌─────────────────────┐      │
│     Disk Device     │◄─────┘
└─────────────────────┘
```

# Isolating components
## discrete elements of an I/O request

- Bottlenecks can happen at any of the layers in either direction

- Isolate the I/O test at one layer
  - mincache = direct
  - /dev/vg/rlvol
  - /dev/rdsk/cxtxdx

- Only test reading or writing

# Isolating components
## discrete elements of an I/O request

**Code Fragment for Timing IO Requests**

```c
#include <sys/time.h>
#define delta_tv(tv_0, tv_1) \
    (tv_1.tv_sec - tv_0.tv_sec + (tv_1.tv_usec - tv_0.tv_usec)/1000000.0)

struct timeval xtv0, xtv1;
struct timezone tz;
double rdt = 0.0;


main()
{
    gettimeofday(&xtv0, &tz);
    read(fd,buf,bufsize)
    gettimeofday(&xtv1, &tz);

    rdt = delta_tv(xtv0, xtv1);

    printf("milliseconds for read:%.3lf ms)\n", 1000*rdt);
}
```

# Isolating components
## taking the network out of the picture

- Multi-tiered applications (e.g. SAP) have large network components which can have a large impact on overall throughput

- Database access is often through sockets

- Techniques for isolation
  - Make local queries rather than client queries
  - With system issues execute problematic commands at the console
  - Use programs/benchmarks similar to those used for IO testing

# Isolating components
## OmniBack debugging techniques

- Understand the capabilities of each component in the configuration

- Isolate
  - Disk I/O
  - Network
  - Tape I/O
  - Updates to OmniBack database
  - Data compressibility

# Knowing your system

- Transaction reporting
  - Example:    SAP instrumentation
    ARM instrumentation

- Maintain a history
  - sar, vmstat, scope, application measures

- Develop an intuition for your systems

- Watch it closely when its healthy

- Know the performance pattern over the day/week/month

- Internals knowledge of the application/database

- Internals knowledge of the OS

# Rules of thumb

- CPU

- Memory

- I/O

# Rules of thumb
## CPU

- System CPU <= 30%

- Total CPU < 80%

- Small load average

# Rules of thumb
## Memory

- Never page out

- Never deactivate processes

# Rules of thumb
## IO

- Utilization < 50% on any drive

- Minimal queuing < 4

- Response time ~10 milliseconds

Interex, Encompass and HP bring you a powerful new HP World.