# Analyzing NFS Client Performance with IOzone

1

Don Capps

Performance Architect

HP

capps@iozone.org

Tom McNeal

Independent Consultant

TMCN Consulting

trmcneal@attbi.com

# Benchmark Overview

# Characteristics of IOzone Activities

# Load Generation

- ## File System I/O requests
  - File sizes vary from 64K to 512M
  - Record sizes vary from 4K to 16M
  - Each increase doubles previous size
  - Large file system calls supported

- ## System variants supported
  - Memory mapped files
  - fread(), fwrite()
  - pread(), pwrite()

# Sequential Reads/Writes

- **Reads & Rereads**
- **Writes & Rewrites**
- **Backwards sequential read**
- **"Stride" read**
  - *Uses constant intervals for sequential reads from beginning to end*

# Other Reads/Writes

- Randomized Reads/Writes

- Record Rewrite (from offset 0)

- fread() – Reads and Rereads
  - Serialized, Buffered & Blocked IO

- fwrite() – Writes and Rewrites
  - Serialized, Buffered & Blocked IO

# Recommended Variants for NFS Clients

./iozone  –azc  –U  /mnt/testdir  –f  /mnt/testdir/testfile

- All tests, all record sizes

- Commit time included in measurements

- IO targeted at mounted file

  – Unmount clears out caches between tests

  – Target file specified in mounted directory

# Characteristics of IOzone Reports

# Graphical Reports
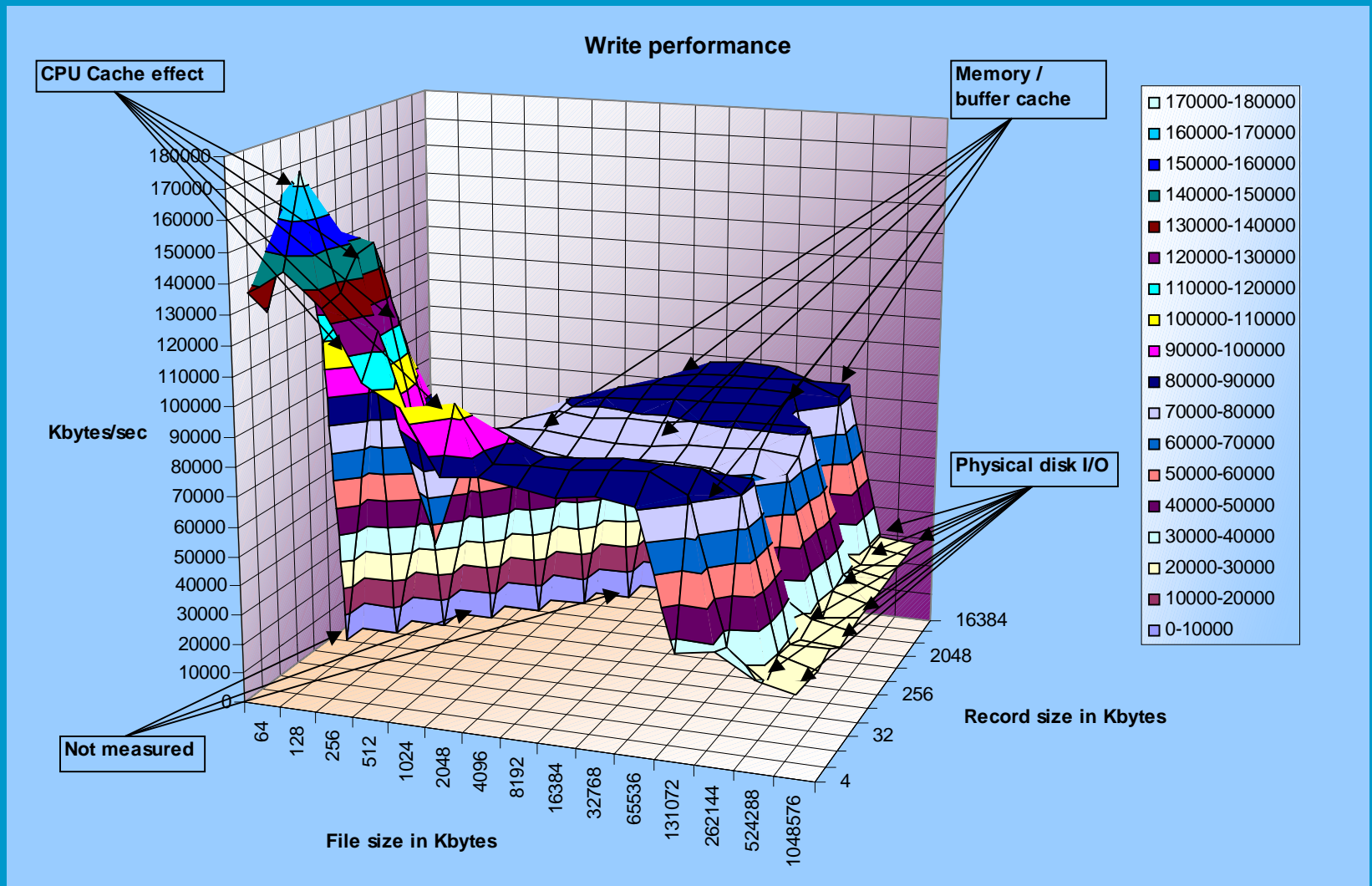
./iozone  –R  -b  exceloutput.xls  > logfile

- Generate Excel output text

- Named file has graphs and data

  – 3D Surface Charts for all tests

  – Includes text output used for graphs

- Standard output sent to log file

  – Generally useful for debugging problems

# Surface Plot Graphs



Write performance

CPU Cache effect

Memory / buffer cache

Physical disk I/O

Not measured

Kbytes/sec

Record size in Kbytes

File size in Kbytes

Legend:
- 170000-180000
- 160000-170000
- 150000-160000
- 140000-150000
- 130000-140000
- 120000-130000
- 110000-120000
- 100000-110000
- 90000-100000
- 80000-90000
- 70000-80000
- 60000-70000
- 50000-60000
- 40000-50000
- 30000-40000
- 20000-30000
- 10000-20000
- 0-10000

# Surface Plot Graphs II

# System Level Variations

*-p, -P #, -I n, -S size, -L size*

- SMP Issues

  – Processor cache purges

  – Processor affinity (for a given # of cpus)

  – Lower bound of number of cpus

- Cache Management

  – CPU Cache size

  – CPU Cache line size

# System parameters

- Client BIOD Daemons

- Server NFSD Daemons

- Number of file system nodes

    – rnode/inode/vnode/file handles

- Directory Name Lookup Cache

- Network buffer sizes

# File System Variations

-o,   -W,   -e,   -g #

- O_SYNC file option for all tests
- File locking required for all IO
- Flush timings included
  - fsync() and fflush()
- Large file offsets
  - File system calls determined at make time
  - Alternate max file size may be specified

# File System Variations II

-B,   -D,   -G,   -H n,   -k n

- Memory mapped file IO
    - mmap() interface
    - MS_ASYNC or MS_SYNC usage available
- Posix asynchronous IO

# Network Variations

- UDP/TCP Protocol

- Client transfer sizes

- Network speed, duplex settings

  – Autonegotiation is often "interesting"

- IP issues

  – Jumbo frames with gigabit ethernet

  – Stream heads, Socket buffer sizes

# Clustered Clients

# Managing and Measuring a Cluster

# Client Specification

-+m  filename

- Clients specified in a file

- Clients must be accessible

  - Remote shells enabled through .rhosts

- DNS $^®$

- IOzone revision 3.128 or later

- Stonewalling helpful (removed by –x)

# "Stonewalling"

- ## Client tests initiated in tandem

  - All clients kept equally busy

- ## When one finishes, they all finish

  - Tests halted when the first client completes

- ## Emulates high performance parallel processing clusters

  - Beowulf clusters at LLNL, PNNL, Los Alamos

# Summary

Examples

and

References

**NFS Industry Conference**

# Summary

- NFS Client measurement standard

  ./iozone –azcR –U /mnt/testdir –f /mnt/testdir/testfile \
  -b exceloutput.xls > logfile

- Gather standard data first

  – What is right for your client?

- Review Variations and Features

  – Review System, FS, and Network setup

  – Start tuning, playing, tuning, playing….

# Rewrite Graph



Re-write performance

**CPU Cache effect**

**Memory / buffer cache**

**Physical disk I/O**

**Not measured**

Kbytes/sec

File size in Kbytes

Record size in Kbytes

Legend:
- 420000-450000
- 390000-420000
- 360000-390000
- 330000-360000
- 300000-330000
- 270000-300000
- 240000-270000
- 210000-240000
- 180000-210000
- 150000-180000
- 120000-150000
- 90000-120000
- 60000-90000
- 30000-60000
- 0-30000

# Reread Graph



Re-read performance

CPU Cache effect

Memory / buffer cache

Physical disk I/O

Not measured

Kbytes/sec

Record size in Kbytes

File size in Kbytes

Legend:
- 700000-750000
- 650000-700000
- 600000-650000
- 550000-600000
- 500000-550000
- 450000-500000
- 400000-450000
- 350000-400000
- 300000-350000
- 250000-300000
- 200000-250000
- 150000-200000
- 100000-150000
- 50000-100000
- 0-50000

# Random Read Graph



**Random Read Performance**

CPU Cache effect

Memory / buffer cache

Physical disk I/O

Not measured

Kbytes/sec

File size in Kbytes

Record size in Kbytes

Legend:
- 560000-600000
- 520000-560000
- 480000-520000
- 440000-480000
- 400000-440000
- 360000-400000
- 320000-360000
- 280000-320000
- 240000-280000
- 200000-240000
- 160000-200000
- 120000-160000
- 80000-120000
- 40000-80000
- 0-40000

# References

- http://www.iozone.org

- http://www.iozone.org/src/current

  - Contains 8K vs. 32K Transfer Size graphs

- http://www.sourceforge.net/projects/nfstestmatrix

  - Includes functional tests, destructive tests, and benchmarks for Linux systems

- http://www.mclx.com/orph

  - Linux Server performance review (late 2001)