

今までのファイルシステムとは一味違う ZFS を活用して面倒なファイル管理の苦勞を減らしましょう！

パフォーマンス改善を行う時、負荷に対して処理能力を持て余している部分を補強してもほとんど改善は望めません。一方、処理能力を目一杯使っている部分を補強すると、大きな改善を期待することができます。

前回は、システム情報を収集するためのツールとして DTrace の基本を紹介させていただきました。今回は、いくつかのサンプルプログラムを使用したシステム情報の収集方法について紹介します。

最初はどのプロセスがどのデバイスを使用しているのかを表示するプログラムです (Solaris 動的トレースガイド P.319)

```
#!/usr/sbin/dtrace -s

#pragma D option quiet

io:::start
{
    @[args[1]->dev_statname, execname, pid] = sum(args[0]->b_bcount);
}

END
{
    printf("%10s %20s %10s %15s\n", "DEVICE", "APP", "PID", "BYTES");
    printa("%10s %20s %10d %15@d\n", @);
}
```

テキストエディタでこのプログラムが書かれたファイルを作成し、以下のように実行してください (ファイル名を ioproc.d) とします。

```
# chmod 700 ioproc.d
# ./ioproc.d
```

```
^C
      DEVICE                APP          PID          BYTES
      sd1                    sched         0             65536
      sd1                    zpool-syspool 5            16625152
      zfs2                   pageout       2            17784832
```

自動的に終了しないので Ctrl + C を押して終了させます。DEVICE は I/O デバイス名、APP はプロセスの実行名、BYTES は入出力が行われたバイト数です。

以下のプログラムはデバイス毎の I/O 速度を表示します (Solaris 動的トレースガイド P.320)

```
#!/usr/sbin/dtrace -s

#pragma D option quiet

io:::start
{
    start[args[0]->b_edev, args[0]->b_blkno] = timestamp;
}

io:::done
/start[args[0]->b_edev, args[0]->b_blkno]/
{
    this->elapsed = timestamp - start[args[0]->b_edev, args[0]->b_
blkno];
```

```

                                cn111
    @[args[1]->dev_statname, args[1]->dev_pathname] =
    quantize((args[0]->b_bcount * 976562) / this->elapsed);
    start[args[0]->b_edev, args[0]->b_blkno] = 0;
}

END
{
    printa(" %s (%s)¥n%d¥n", @);
}

```

テキストエディタでこのプログラムが書かれたファイルを作成し、以下のように実行してください（ファイル名を iospeed.d）とします。

```

# chmod 700 iospeed.d
# ./iospeed.d

```

```

dtrace: aggregation size lowered to 1m
^C

```

```

sd1 (/devices/pci@0,0/pci1179,1@1d/storage@2/disk@0,0:r)
value  ----- Distribution ----- count
   32 |                                     0
   64 |                                     3
  128 |                                     1
  256 | @@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@ 2257
   512 |                                     1
  1024 |                                     0

```

```

sd2 (/devices/pci@0,0/pci8086,3410@9/pci1000,9280@0/sd@1,0:a)
value  ----- Distribution ----- count
  128 |                                     0
  256 |                                     1
  512 |                                     0
 1024 |                                     2
 2048 |                                     0
 4096 |                                     2
 8192 | @@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@ 172
16384 | @@@@@@                               52
32768 | @@@@@@@@@@@@@@@@@@                 108
65536 | @@@@                                34
131072|                                     0

```

実行すると計測を開始します。そして Ctrl + C が押されるまで情報を収集し、デバイス毎の統計情報を表示します。iostat でも同じ情報を収集することができますが、このプログラムの方がより分かりやすい形で結果が表示されます。sd1 と sd2 を比較すると sd1 は 256 ~ 512 [KB/sec]、sd2 は 8 ~ 32 [MB/sec] の速度で読み書きが行われていることが分かります。

例えば、これが sd1 から読み出しして sd2 へ書くような処理だった場合、sd1 が律速であることが分かります。このような処理のパフォーマンスを改善するには sd1 を早いデバイスにすると効果的です。

次のプログラムは write() システムコールが復帰するまでの時間を、プロセス毎に表示するプログラムです。

```

#!/usr/sbin/dtrace -s

syscall::write:entry
{
    self->ts = timestamp;
}

```

cnl11

```
syscall::write:return
/self->ts/
{
    @time[execname] = quantize(timestamp - self->ts);
    self->ts = 0;
}
```

テキストエディタでこのプログラムが書かれたファイルを作成し、以下のよう
に実行してください（ファイル名を writetm.d）とします。

```
# chmod 700 writetm.d
# ./writetm.d
```

```
dtrace: script './writetm.d' matched 2 probes
dtrace: aggregation size lowered to 2m
^C
```

```
tar
value ----- Distribution ----- count
  512 | 0
 1024 | 20
 2048 | @@@@ 23505
 4096 | @@@@ 31122
 8192 | @@@@ 51358
16384 | @@@@ 34861
32768 | @@ 7059
65536 | 1531
131072 | 381
262144 | 135
524288 | 39
1048576 | 3
2097152 | 17
4194304 | 46
8388608 | 27
16777216 | 36
33554432 | 50
67108864 | 40
134217728 | 6
268435456 | 0
```

実行すると計測を開始します。そして Ctrl + C が押されるまで情報を収集し、プロセス毎の統計情報を表示します。ここでは tar についての結果のみを抜粋しています。

プログラムの write の部分を read に変更すれば、read() システムコールの所要時間を得ることもできます。timestamp は D の組み込み変数でナノ秒単位のカウンタです。このプログラムのように相対的な時間を計測する場合に使用することができます。

今回のコラムでは、DTrace を使用したシステム情報の収集例を紹介させていただきました。次回は ARC の情報収集について紹介します。

(つづく)

■ 参考文献

- (1) サンマイクロシステムズ、“SunOS リファレンスマニュアル 1 : ユーザコマンド”、Part No. 819-1210-13、2007 年 7 月
- (2) サンマイクロシステムズ、“SunOS リファレンスマニュアル 1M : システム管理コマンド”、Part No. 819-1211-13、2007 年 7 月
- (3) サンマイクロシステムズ、“Solaris 動的トレースガイド”、Part No. 819-0395-12、2008 年 10 月

