# NFS over TCP:

## Excessive wakeups with BSD-based stacks

Ric Werme

Compaq Computer

Connectathon 2002

# A single 64 KB NFS write

- 45 "classic" frames or
- 8 jumbo frames:

clnt.791 > svr.2049: P 829:9789(8960) ack 873 win 49152

clnt.791 > svr.2049: P 9789:18749(8960) ack 873 win 49152

clnt.791 > svr.2049: P 18749:27709(8960) ack 873 win 49152

clnt.791 > svr.2049: P 27709:36669(8960) ack 873 win 49152

clnt.791 > svr.2049: P 36669:45629(8960) ack 873 win 49152

clnt.791 > svr.2049: P 45629:54589(8960) ack 873 win 49152

clnt.791 > svr.2049: P 54589:63549(8960) ack 873 win 49152

clnt.791 > svr.2049: P 63549:66513(2964) ack 873 win 49152

svr.2049 > clnt.791: . ack 18749 win 40192

svr.2049 > clnt.791: . ack 36669 win 31232

svr.2049 > clnt.791: . ack 54589 win 22272

svr.2049 > clnt.791: . ack 66513 win 49152

svr.2049 > clnt.791: P 873:1037(164) ack 66513 win 49152

# Data throttles

- TCP throttles:

  window

  congestion control

- Sockbuf throttle:

  sb_hiwat/soreserve()

- Atomic sends

  Can't interleave RPC messages!

# Data release

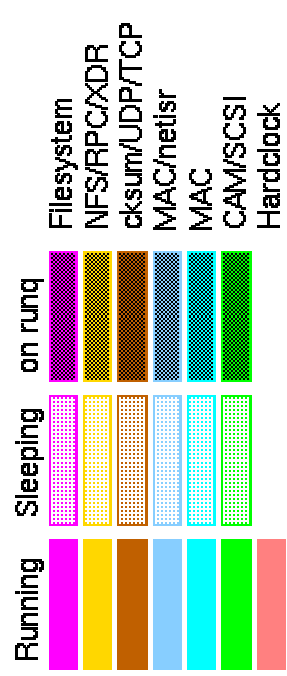- TCP window updates
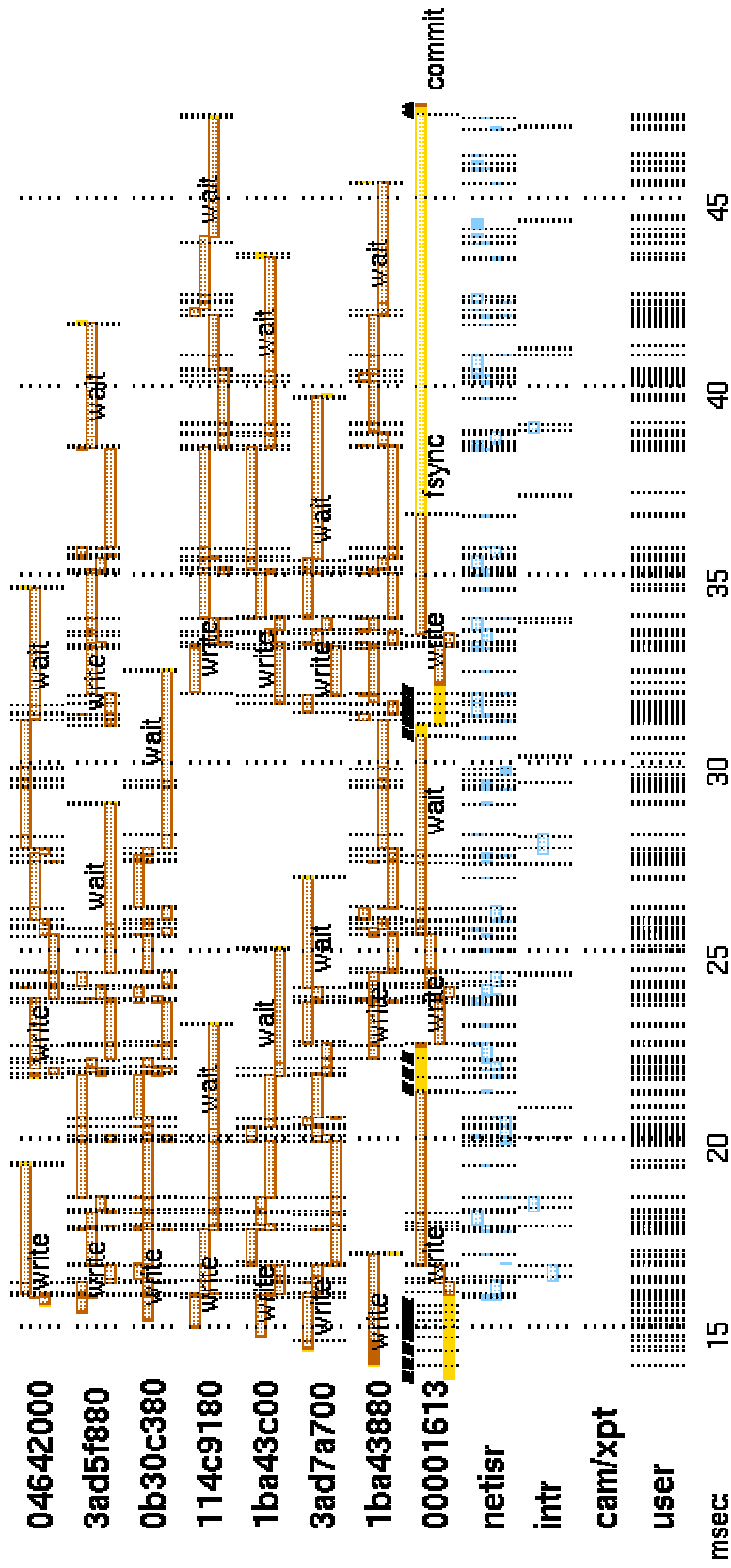- TCP ACKs and sockbuf removals

  svr.2049 > clnt.791: . ack 54589 win 22272

  svr.2049 > clnt.791: . ack 66513 win 49152

# Multiple NFS writes

- Client: Sends several more 64 KB writes. Eventually, nfsiod threads block waiting for 64 KB in the sockbuf.

- As ACKs arrive, acknowledged data is removed and sowwakeup() called to let "the" sender buffer more data.

- All the client threads wakeup, typically all find nothing has changed and go back to sleep.

- Sometimes one thread gets lucky and buffers the next write.
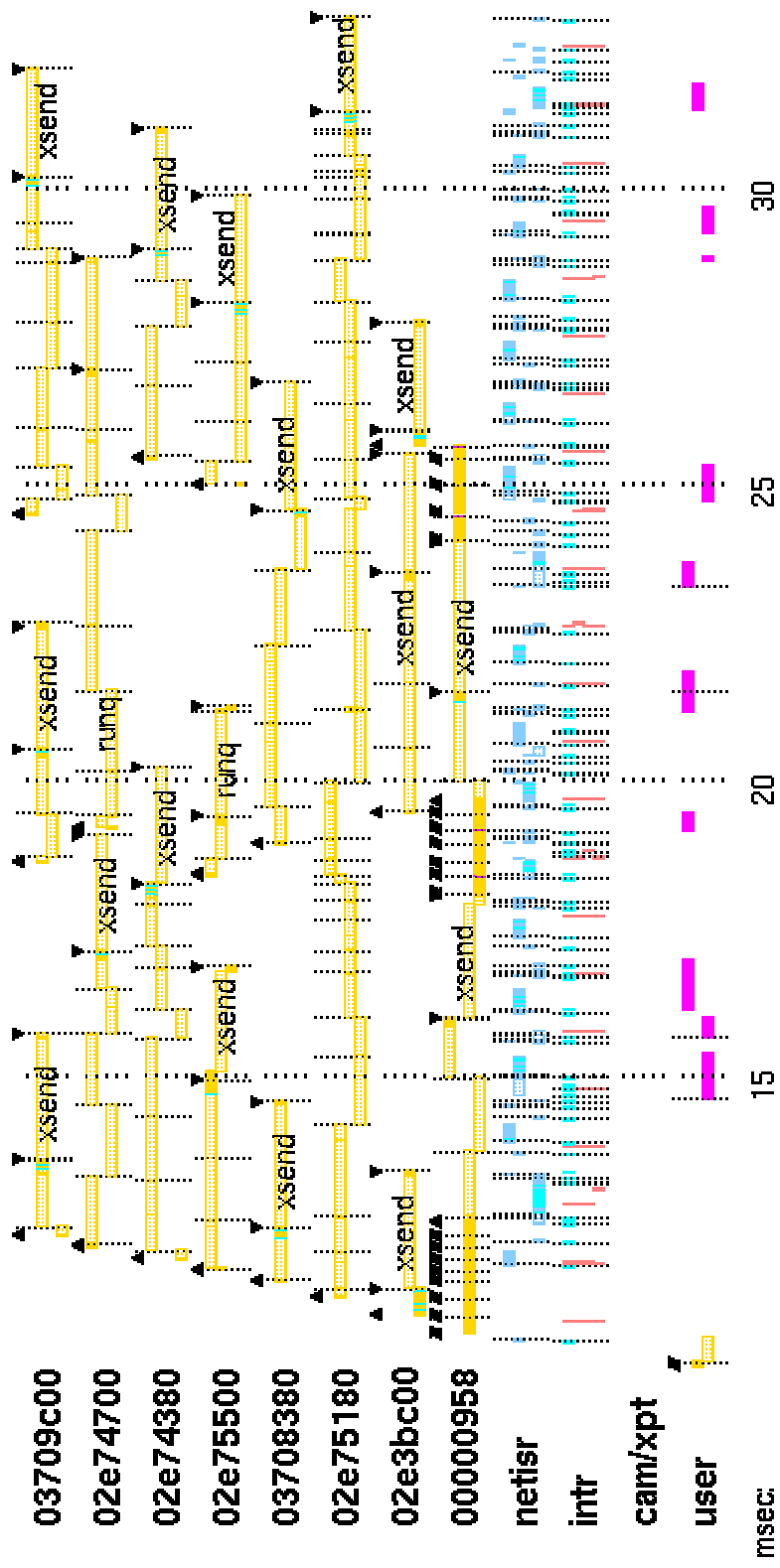
# 16 64 KB writes, 96 KB window



04642000
3ad5f880
0b30c380
114c9180
1ba43c00
3ad7a700
1ba43880
00001613
netisr
intr
cam/xpt
user

msec:    15        20        25        30        35        40        45

Running    Sleeping    on runq

Filesystem
NFS/RPC/XDR
cksum/UDP/TCP
MAC/netisr
MAC
CAM/SCSI
Hardclock

Input: iod7.96k
Output: iod7.fig
Start time: 13.642 msec
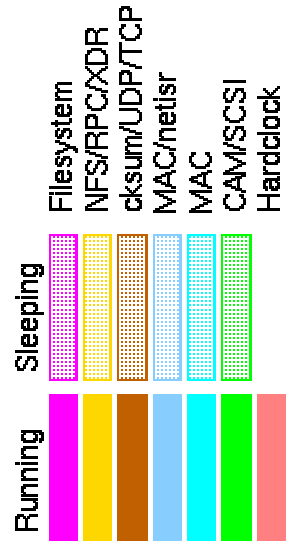Date: Fri Feb 15 13:49:12 2002
By: werne

# SB_WAKEONE?

- BSD can awaken a single thread in sowakeup(), currently not used in Tru64 Unix. (Other mechanisms are used in NFS and elsewhere.

- Does not solve the reordering problem as requests remain "SIRO" – sequential in, random out.

# 64 KB TCP client writes to 1152 KB file
## 98304 byte TCP window, 7 nfsiod threads, SB_WAKEONE

03709c00
02e74700
02e74380
02e75500
03708380
02e75180
02e3bc00
00000958
netisr
intr
cam/xpt
user

msec:

xsend
xsend
xsend
xsend
xsend
xsend
xsend
xsend
xsend
xsend
xsend
xsend
runq
runq

15
20
25
30

**Running**
**Sleeping**

Filesystem
NFS/RPC/XDR
cksum/UDP/TCP
MAC/netisr
MAC
CAM/SCSI
Hardclock

ES40 (four CPU) client and server

# Other solutions

- ## Queue of threads waiting for sockbuf

  Extra locking, code, overhead.

  Could be done mostly in RPC code.

  Could use sb_wakeup() callback.


- ## Bigger sockbuf

  Trivial change, at least for a first pass.

  Memory is cheap.  Really cheap!

  Bigger TCP window on receive side,

  permits bigger bandwidth-delay product.

# 16 64 KB writes, 512 KB window



0b30c380
04642000
3ad7a700
1ba43880
1ba43c00
3ad5f880
114c9180
00003518
netisr
intr
cam/xpt
user

msec:
15    20    25    30    35    40    45

write    commit

Running    Sleeping    on runq

Filesystem
NFS/RPC/XDR
cksum/UDP/TCP
MAC/netisr
MAC
CAM/SCSI
Hardclock

Input: iod7.512k
Output: iod7.fig
Start time: 11.594 msec
Date: Fri Feb 15 13:47:44 2002
By: werme

# NFS reads

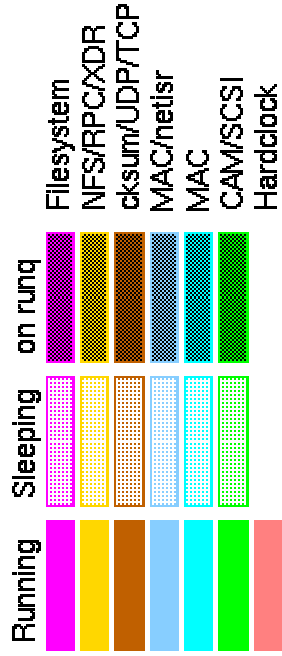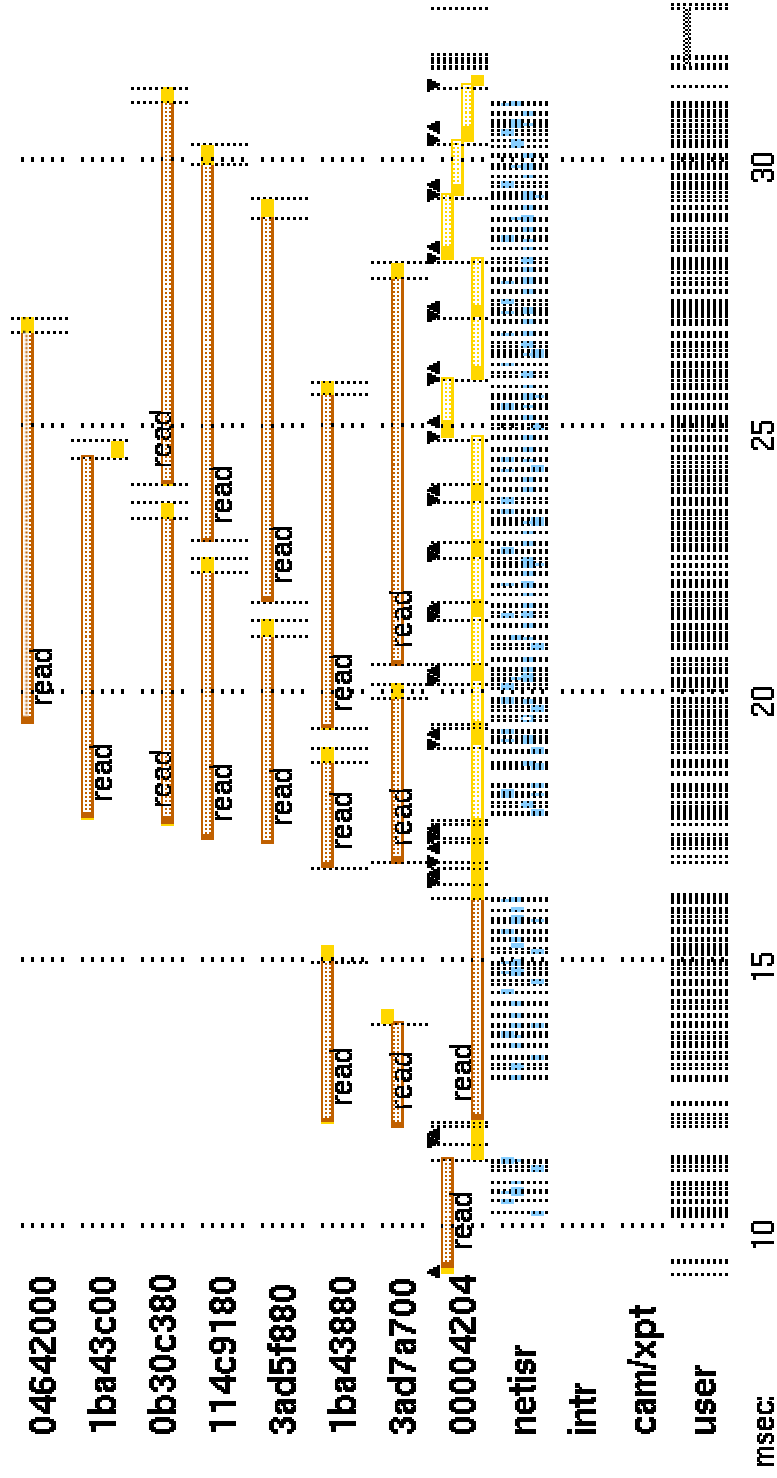- Readahead done differently than write behind

  Block 0: no read ahead

  Block 1: issue read aheads for 2 & 3, read 1

  Block 2: issue read aheads for 4 & 5, wait for 2
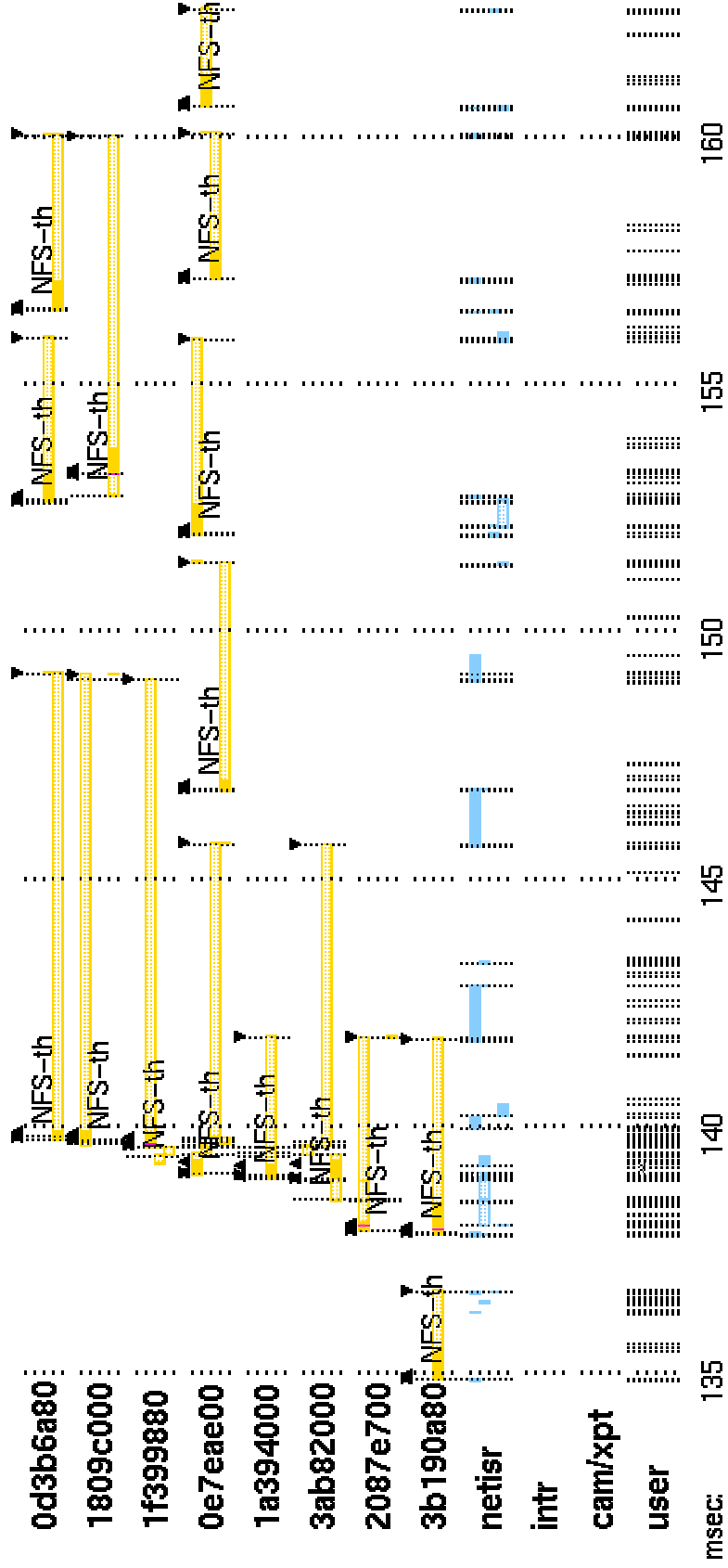
  Block n: issue read ahead for n+8, wait for n

- First, client trace (shows reads better)

- Second, server trace

- No big surprises

# 16 64 KB reads, 512 KB window, client view



04642000
1ba43c00
0b30c380
114c9180
3ad5f880
1ba43880
3ad7a700
00004204
netisr
intr
cam/xpt
user

msec:          10          15          20          25          30

| Running | Sleeping | on runq | |
|---------|----------|---------|---|
| | | | Filesystem |
| | | | NFS/RPC/XDR |
| | | | cksum/UDP/TCP |
| | | | MAC/netisr |
| | | | MAC |
| | | | CAM/SCSI |
| | | | Hardclock |

Input: cfkpread512.kfig
Output: cfkpread512.fig
Start time: 9.083 msec
Date: Tue Feb 19 14:09:36 2002
By: weme

# 16 64 KB reads, 512 KB window, server view



0d3b6a80  
1809c000  
1f399880  
0e7eae00  
1a394000  
3ab82000  
2087e700  
3b190a80  
netisr  
intr  
cam/xpt  
user  

msec: 135 140 145 150 155 160

Running    Sleeping    on runq

Filesystem  
NFS/RPC/XDR  
cksum/UDP/TCP  
MAC/netisr  
MAC  
CAM/SCSI  
Hardclock

Input: suread512.ktg  
output: suread512.fig  
Start time: 134.848 msec  
Date: Tue Feb 19 14:19:49 2002  
By: werne

# 16 64 KB reads, 96 KB window

28254e00
21bf4000
0d3b7180
3ae01880
23355c00
1f8c3500
3a83b880
27900000
netisr
intr
cam/xpt
user

NFS-th
NFS-th
NFS-th
NFS-th
NFS-th
NFS-th
NFS-th
NFS-th

NFS-th
NFS-th
NFS-th
NFS-th
NFS-th
NFS-th
NFS-th
NFS-th

msec:    180        185        190        195        200        205        210

Running          Sleeping          on runq

Filesystem
NFS/RPC/XDR
cksum/UDP/TCP
MAC/netisr
MAC
CAM/SCSI
Hardclock