

TCP/IP performance test

y.ma@riken.jp

20170407

Outline

- ❖ Motivation:

- ❖ Measure the real performance and prepared as input for more integrated study

- ❖ Setup: (more details in slides)

- ❖ Software: mini-DAQ
- ❖ Hardware: 2 sets of 40 core machine with direct 10gbps connection

- ❖ Results

- ❖ Summary & todo

Setup: Hardware

- ❖ Test performance with “*iperf*” tools
- ❖ Test performance with “*scp*”: 1.6gbps, limited by hard drive throughput?

```
meson ~ $ e50_server1_login
oper@192.168.10.202's password:
Last login: Wed Apr  5 15:17:25 2017 from gateway
[oper@e50_server1 ~]$ iperf -s
-----
Server listening on TCP port 5001
TCP window size: 85.3 KByte (default)
-----
[ 4] local 10.0.0.13 port 5001 connected with 10.0.0.3 port 43422
-----
Client connecting to 10.0.0.3, TCP port 5001
TCP window size: 1.05 MByte (default)
-----
[ 6] local 10.0.0.13 port 48888 connected with 10.0.0.3 port 5001
Waiting for server threads to complete. Interrupt again to force quit.
[ ID] Interval      Transfer    Bandwidth
[ 6] 0.0-10.0 sec  10.8 GBytes  9.32 Gbits/sec
[ 4] 0.0-10.0 sec  10.9 GBytes  9.34 Gbits/sec
[oper@e50_server1 ~]$
```

```
meson ~ $ e50_server0_login
oper@192.168.10.201's password:
Last login: Fri Apr  7 12:30:40 2017 from gateway
[oper@e50_server0 ~]$ iperf -c 10.0.0.13 -d
-----
Server listening on TCP port 5001
TCP window size: 85.3 KByte (default)
-----
Client connecting to 10.0.0.13, TCP port 5001
TCP window size: 1.14 MByte (default)
-----
[ 5] local 10.0.0.3 port 43422 connected with 10.0.0.13 port 5001
[ 4] local 10.0.0.3 port 5001 connected with 10.0.0.13 port 48888
[ ID] Interval      Transfer    Bandwidth
[ 5] 0.0-10.0 sec  10.9 GBytes  9.35 Gbits/sec
[ 4] 0.0-10.0 sec  10.8 GBytes  9.31 Gbits/sec
[oper@e50_server0 ~]$
```

Setup: Software

- ❖ Default CentOS 7.0 x86_64:

- ❖ [oper@e50_server0 ~]\$ cat /proc/sys/net/ipv4/tcp_rmem

- ❖ 4096 (min [Byte]) 87380(default [Byte]) 6291456(max [Byte])

- ❖ [oper@e50_server0 ~]\$ cat /proc/sys/net/ipv4/tcp_wmem

- ❖ 4096(min [Byte]) 16384(default [Byte]) 4194304(max [Byte])

- ❖ mini-DAQ:

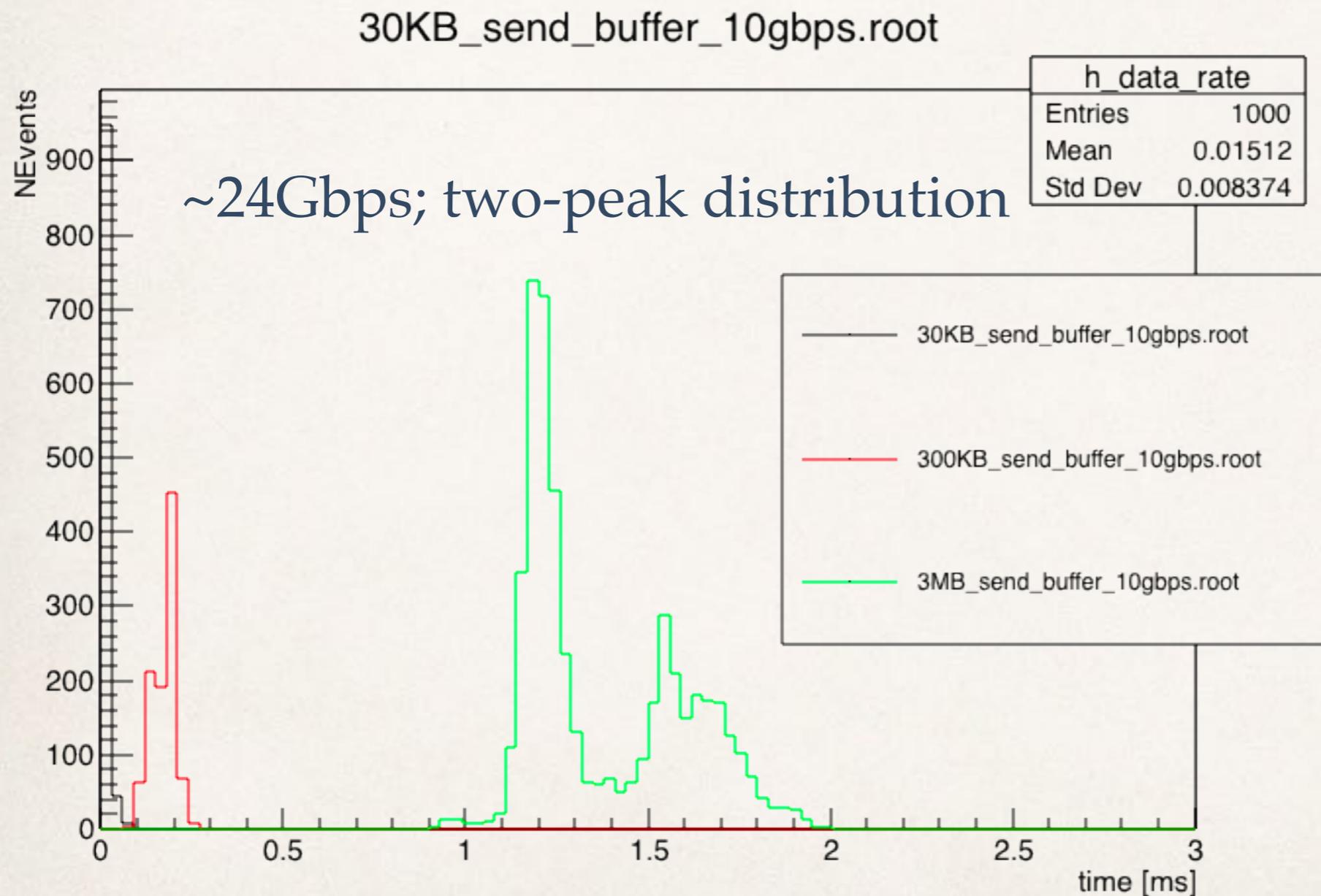
- ❖ minimum DAQ functionality from scratch in C

- ❖ data source (ring buffer) —> *TCP/IP socket* —> data sink (ring buffer)



default C library function,
no customization

Results: data source *send()* throughput

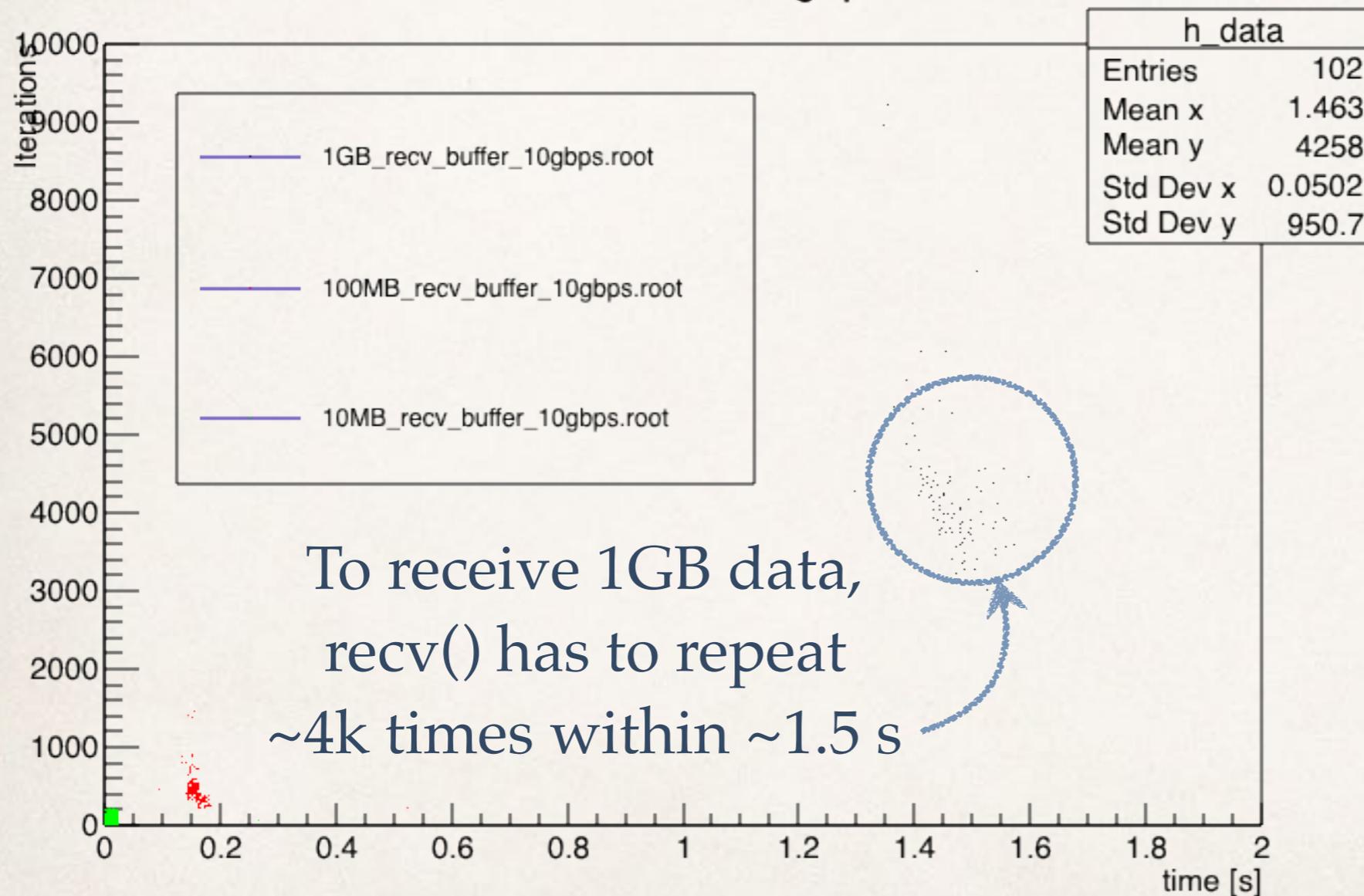


Method:

At data source,
change `send()`
buffer size
and measure the
time consumption;
Up to 30MB,
the `send()` function
works at 24Gbps;
Larger than 30MB,
process hangup...

Results: data sink *recv()* throughput

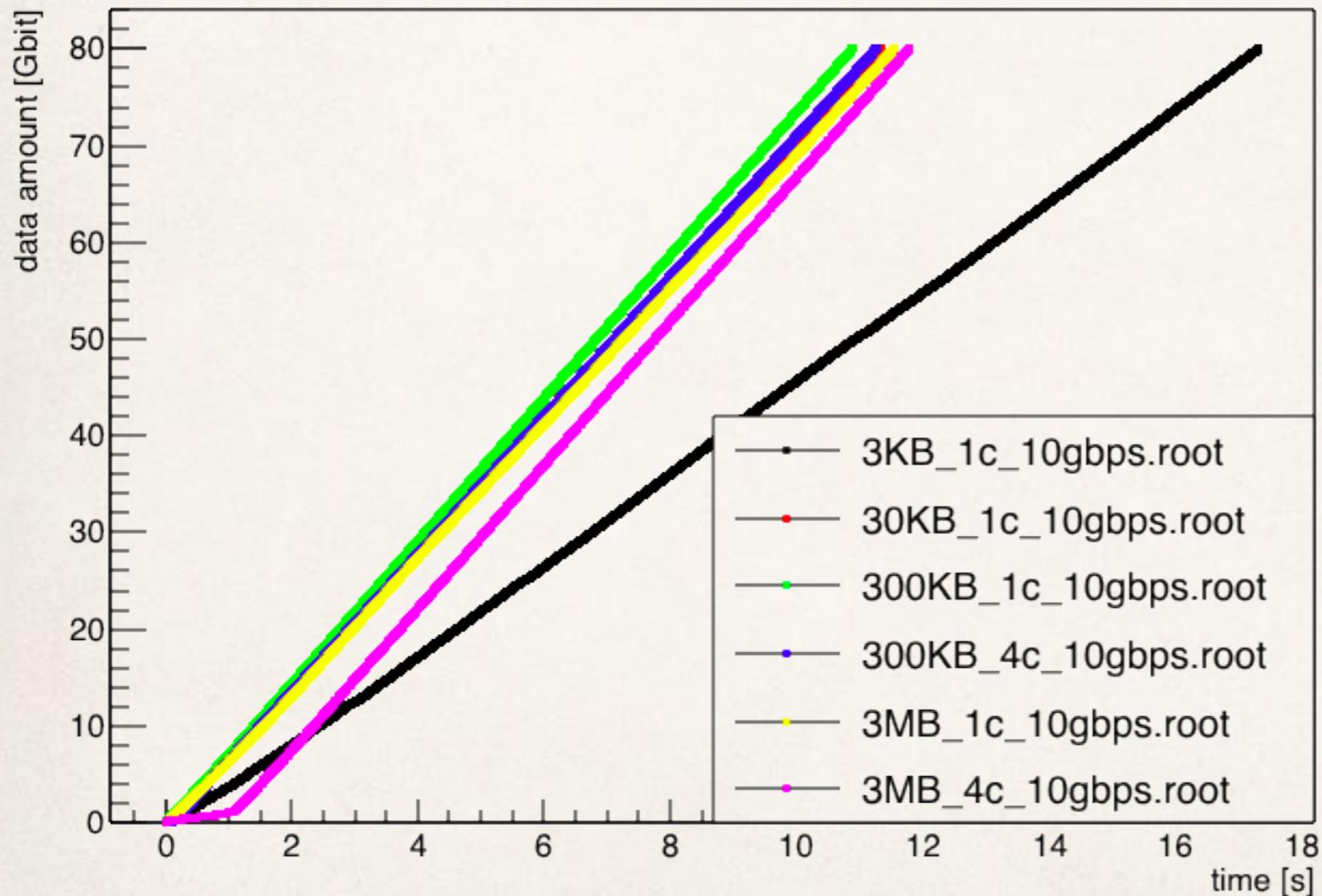
1GB_recv_buffer_10gbps.root



Method:
At data sink,
change `recv()`
buffer size
and measure the
receive cycles and
time consumption;

Results: data *transportation* speed

Ethernet throughput test



Method:

Adjust **buffer size**
on both `recv()` and `send()`,
measure data transfer speed;
Too small buffer size (3KB)
results in low performance;
multi clients (data source)
have similar speed,

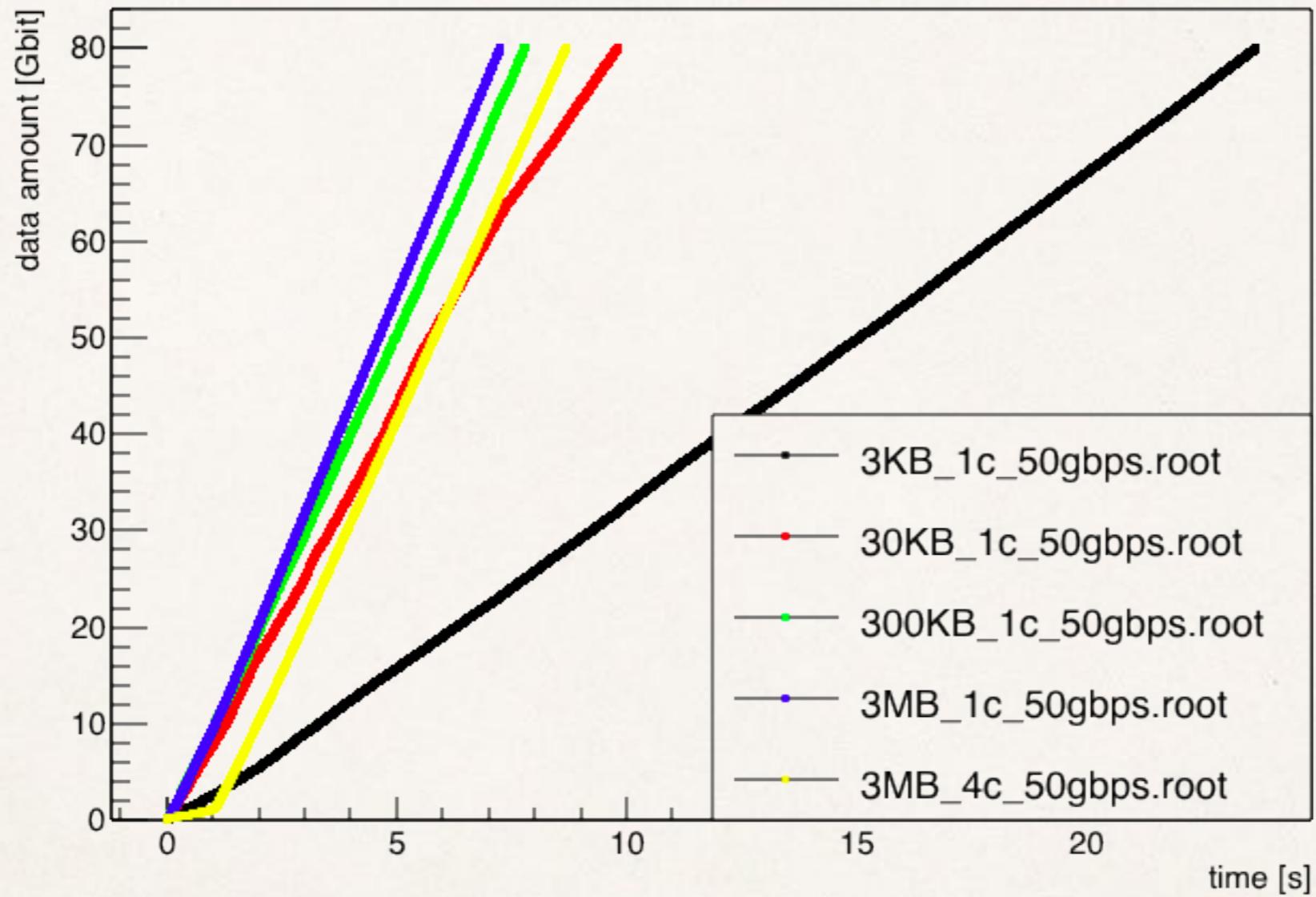
*though it takes
some time to
reach the optimized point.*

Summary & todo

- ❖ Single, direct link 10Gbps ethernet has reasonable performance;
- ❖ The jittering is not negligible and a *safety factor of two* seems necessary for a reliable performance estimation
- ❖ Use 10Gbps switch to test data collision effect(Cisco UCS 6120xp);
- ❖ At the moment, 50Gbps seems have a performance bottle neck issue...

❖ Backup

Ethernet throughput test

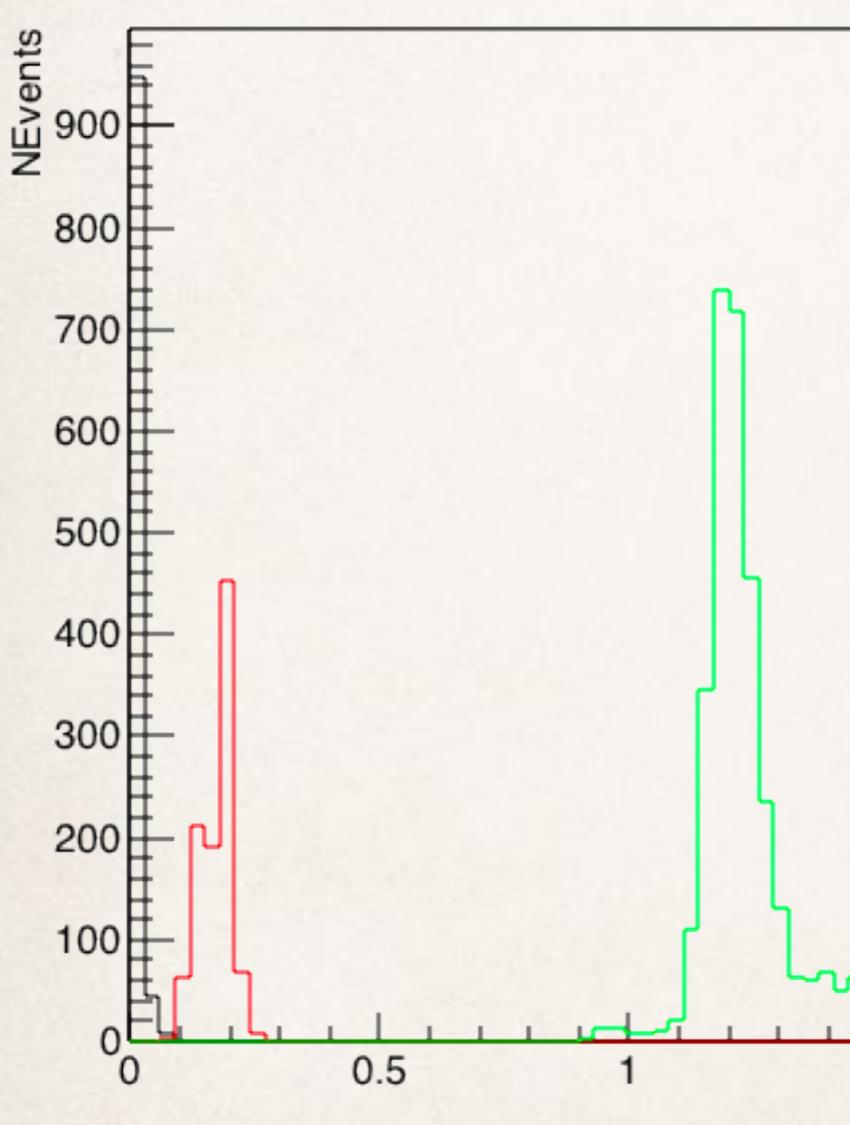


```

int send_all_tuning(int sockfd, void *buff, uint32_t len, double *time_diff, uint32_t *counter)
{
    uint32_t total = 0;
    uint32_t bytesleft = len;
    uint32_t n;
    timespec start_time, stop_time;
    *counter = 0;
    clock_gettime(CLOCK_PROCESS_CPUTIME_ID, &start_time);
    while(total < len)
    {
        n = send(sockfd, buff + total, bytesleft, 0);
        if (n == -1)
            break;
        total += n;
        bytesleft -= n;
        (*counter)++;
    }
    len = total;
    clock_gettime(CLOCK_PROCESS_CPUTIME_ID, &stop_time);
    *time_diff = (stop_time.tv_sec - start_time.tv_sec) * 1000. + (stop_time.tv_nsec - start_time.tv_nsec) / 1000000.;
    printf("send_all_tuning(): counter = %d, time_diff = %f[ms], data = %c[Byte]\n", *counter, *time_diff, len);
    return n == -1 ? -1 : 0;
}

```

30KB_send_buffer_10gbps.root



h_data_rate	
Entries	1000
Mean	0.01512
Std Dev	0.008374

```

--- utility.c 12% L139 (C/L Abbrev)
double time_tuning;
uint32_t counter_tuning;
if( send_all_tuning( fd socket data client,
                    (uint8_t *)event_data,
                    3000000,
                    &time_tuning,
                    &counter_tuning ) == -1 )
{
    printf_color("r", "data sockets broken!!\n");
    break;
}
h_data_rate->Fill( time_tuning );
evt_crt++;
--- transporter_thrd.c 86% L104 (C/L Abbrev)

```