Hybrid Start A TCP Slow Start Algorithm

Narayan Pai Final year CSE, NITK Surathkal

narayanspai2000@gmail.com Roll No: 181CO135

Outline

- 1. Congestion Control
- 2. Phases in Congestion Control Algorithm
- 3. Problems with standard slow start
- 4. Hybrid Start Algorithm
- 5. Kernel Code Walkthrough
- 6. Experimenting with the kernel



Congestion Control

- → It is a mechanism within TCP that makes sure congestion does not occur in the network.
- → Enforced with a state variable named Congestion Window (*cwnd*) that limits the number of packets the sender can send into the network before receiving an ACK.

Phases in Congestion Control

Slow Start Phase

Starting with a small value, increment cwnd exponentially to ssthresh

Congestion Avoidance Phase

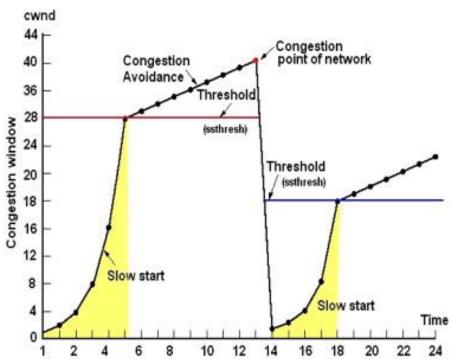
Increment cwnd linearly

Congestion Detection Phase

Decrement cwnd and ssthresh.

^{*}ssthresh (slow start threshold) is a dynamic threshold for cwnd

Phases in Congestion Control Algorithm



Problems with standard slow start

- 1. **Heavy packet losses** occur because of the exponential increase of the congestion window during standard slow start.
- 2. Results in a **long blackout period** without transmission.

For eg, second packet is dropped when cwnd = 100.

Activity: Think about what happens if there are other packet drops too?



Hybrid Start : The solution?

Introduced in 2009 by Sangtae Ha and Injong Rhee in a paper named "Taming the elephants: New TCP slow start"

- → Uses ACK trains and RTT delay samples
- Heuristically finds safe exit points at which it can finish slow start and move to congestion avoidance before cwnd overshoots



Fact

It is used in TCP Cubic, the default TCP in Linux.

Packet train

Set of packets of same length sent back to back from the sender to the receiver to probe the network condition.

If N packets of length L are sent, then the receiver can estimate bandwidth by:

bandwidth by:

$$b(N) = \frac{(N-1)L}{\triangle(N)}$$

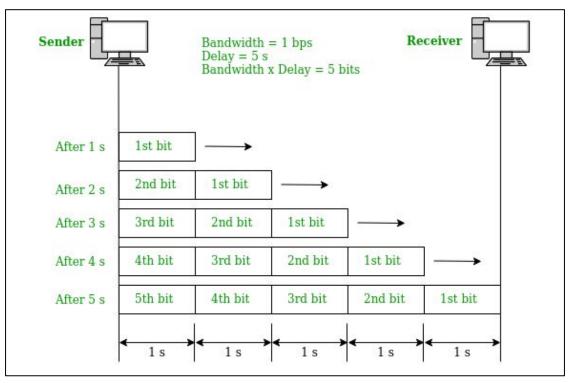
Δ(N) = Packet train length
= Inter-arrival time between
first and the last packet

Bandwidth Delay Product

 Bandwidth delay product is a measurement of how many bits can fill up a network link.

- Gives an upper bound for cwnd

Bandwidth Delay Product



Exit Criterion #1: ACK Train length

Hystart uses the data packets itself as probe packet train.

$$cwnd.L = (N-1).E-(1)$$

During complete network utilization,

$$cwnd.L = b(N) \cdot D_{min} = \frac{(N-1) \cdot L}{\Delta(N)} \cdot D_{min} \qquad --(2)$$

By (1) and (2), we get $\Delta(N) = D_{min}$ during complete network utilization.

Conclusion: Exit when $\Delta(N) > D_{min}$

Implementation Details

- 1. To estimate $\Delta(N)$, we use the train of ACKs received in response to a packet train. Helps in conservative estimation of the available capacity of the path.
- 2. To estimate D_{min} , we divide the minimum RTT by two.

Problems!

- 1. When other flows are competing for the same path, D_{min} doesn't give us the required value.
- Not effective when the network is asymmetric.

Exit Criterion #2: Delay increase

Step 1: Measure minimum RTT of a sample in each round.

Step 2: If it is considerably greater than that of the previous round, exit.

Linux Code Walkthrough

Configuration variables

net/ipv4/tcp cubic.c

```
/* Two methods of hybrid slow start */
#define HYSTART ACK TRAIN 0x1
#define HYSTART DELAY 0x2
/* Number of delay samples for detecting the increase of delay */
#define HYSTART MIN SAMPLES 8
#define HYSTART DELAY MIN (4000U) /* 4 ms */
#define HYSTART DELAY MAX (16000U) /* 16 ms */
#define HYSTART DELAY THRESH(x) clamp(x, HYSTART DELAY MIN, HYSTART DELAY MAX)
static int hystart read mostly = 1;
static int hystart detect read mostly = HYSTART ACK TRAIN | HYSTART DELAY;
static int hystart low window read mostly = 16;
```

```
static void hystart update (struct sock *sk, u32 delay)
      struct tcp sock *tp = tcp sk(sk);
      struct bictcp *ca = inet csk ca(sk);
     u32 threshold;
     if (hystart detect & HYSTART ACK TRAIN) {
         u32 now = bictcp clock us (sk);
          ca->last ack = now;
          threshold = ca->delay min + hystart ack delay (sk);
          if ((s32) (now - ca->round start) > threshold) {
            ca \rightarrow found = 1;
            NET INC STATS (sock net (sk),
                    LINUX MIB TCPHYSTARTTRAINDETECT);
            NET ADD STATS (sock net (sk),
                    LINUX MIB TCPHYSTARTTRAINCWND,
                    tp->snd cwnd);
            tp->snd ssthresh = tp->snd cwnd;
     contd ...
```

```
if (hystart_detect & HYSTART_DELAY) {
   /* obtain the minimum delay of more than sampling packets */
   if (ca->curr_rtt > delay)
    ca->curr rtt = delay;
   if (ca->sample_cnt < HYSTART_MIN_SAMPLES) {</pre>
     ca->sample_cnt++;
   } else {
     if (ca->curr rtt > ca->delay min +
         HYSTART_DELAY_THRESH(ca->delay_min >> 3)) {
       ca -> found = 1;
       NET_INC_STATS(sock_net(sk),
               LINUX_MIB_TCPHYSTARTDELAYDETECT);
       NET_ADD_STATS(sock_net(sk),
               LINUX_MIB_TCPHYSTARTDELAYCWND,
               tp->snd cwnd);
       tp->snd ssthresh = tp->snd cwnd;
```

Experimenting with the kernel

Tools Used: Flent, NeST

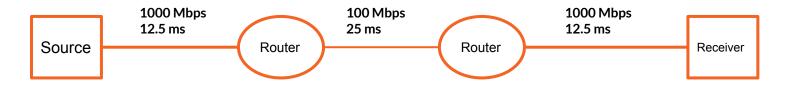
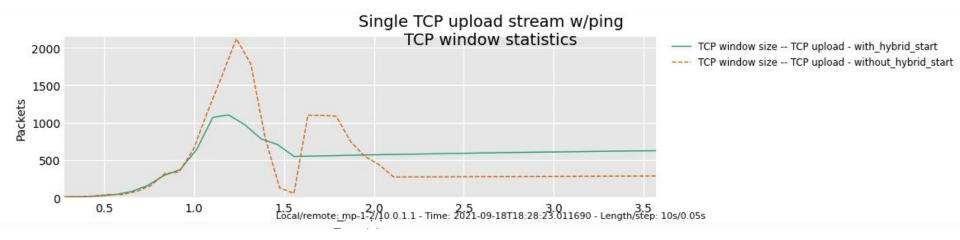


Fig: Experiment Topology with link bandwidth and delay

- Step 1: Run tests with default kernel
- Step 2: Download the kernel source
- Step 3: Customize the kernel
- Step 4: Build and install the kernel
- Step 5: Run tests with the customized kernel

Results



Bibliography

- 1. "Taming the elephants: New TCP slow start": https://doi.org/10.1016/j.comnet.2011.01.014
- 2. Instructions, code and results of the experiment: https://github.com/narayanpai1/hybrid-start-experiments
- 3. Flent tool: https://flent.org
- 4. NeST tool: http://nest.nitk.ac.in

Thank You!