

# Traffic policing in eBPF: applying token bucket algorithm

LPC 2018

Julia Kartseva, [hex@fb.com](mailto:hex@fb.com)

# Traffic policing in eBPF: applying token bucket algorithm

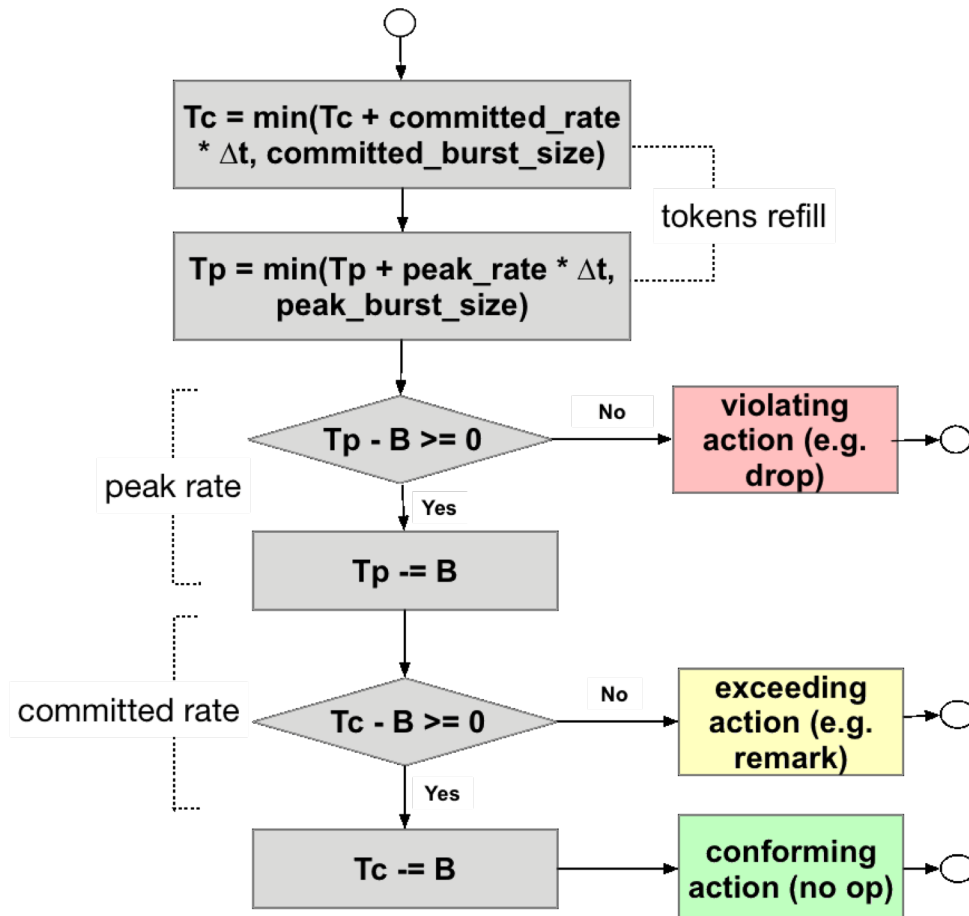
Shaping vs. Policing

Shaping	Policing
Buffers exceeding packets	No buffering, instant action: drop or remark
Latency increase due to buffering	No latency increase, but drops may cause retransmits (e.g. in TCP)
Smooths traffic burstiness, output rate doesn't deviate much	Bursts are propagated
Drops packets anyway when buffer capacity is reached. Buffer increasing causes higher latency	
Linux Traffic Control: queuing disciplines, e.g. tc htb	Switch side policers, eBPF-based traffic policing

# Traffic policing in eBPF: applying token bucket algorithm

RFC 2698, naive implementation

## Two rate three color marker



## eBPF program in TC egress chain

```
1  _u64 delta_t = packet_ts -
2      bucket->timestamp;
3  bucket->tokens += delta_t * rate_bps /
4      NS_IN_SEC;
5  bucket->tokens = MIN(bucket->tokens,
6      burst_size);
7  bucket->timestamp = packet_ts;
8  __u64 tokens_spent = 8 * skb->len;
9  /* TC_ACT_SHOT if no enough tokens */
10 __sync_fetch_and_add(bucket->tokens,
11 (-1) * tokens_spent);
12 return TC_ACT_UNSPEC;
```

Does this code produce the desired rate?

# Traffic policing in eBPF: applying token bucket algorithm

Advancing naive implementation

## Problem:

Updates are in the kernel space. Data race with multi CPU. Getting and adding tokens into a bucket must be executed as an atomic action.

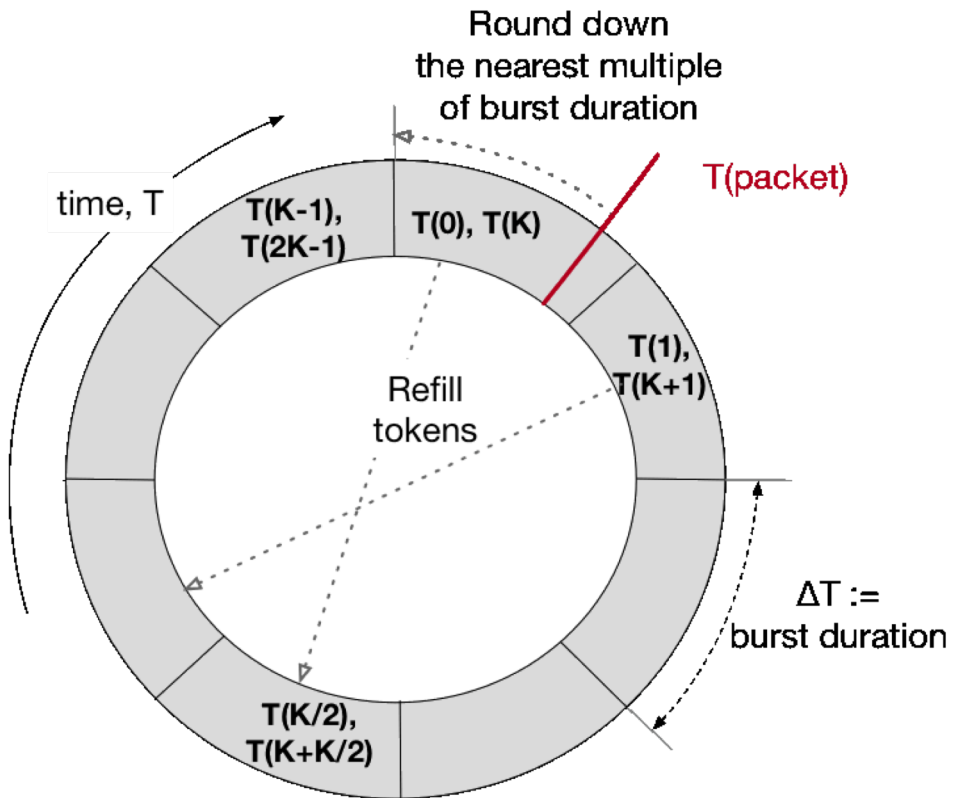
## Solution or not?

- Critical section in eBPF program
- Per CPU eBPF data structures
- Update tokens from the user space  
What if burst duration is in microseconds?
- Data structures shared between CPUs: lru\_hash, array

# Traffic policing in eBPF: applying token bucket algorithm

Hackish working implementation: kernel space only, eBPF array

Key idea: refill tokens in a future bucket, take tokens from the current bucket.

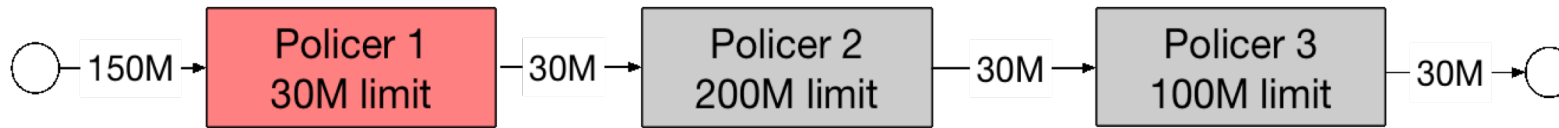


```
1 __u64 packet_ts = bpf_ktime_get_ns();
2 __u64 burst_dur = burst_size * NS_IN_SEC /
3   rate_bps;
4 __u32 refill_tbucket_idx = (packet_ts +
5   (K >> 1) * burst_dur) / burst_dur % K;
6 __s64* refill_tokens = bpf_map_lookup_elem(
7   tbucket_arr, &refill_tbucket_idx);
8 if (refill_tokens) *refill_tokens =
9   burst_size;
10 __u32 curr_tbucket_idx = packet_ts /
11   burst_dur % K;
10 /* Subtract tokens from curr_tbucket_idx
11   bucket */
```

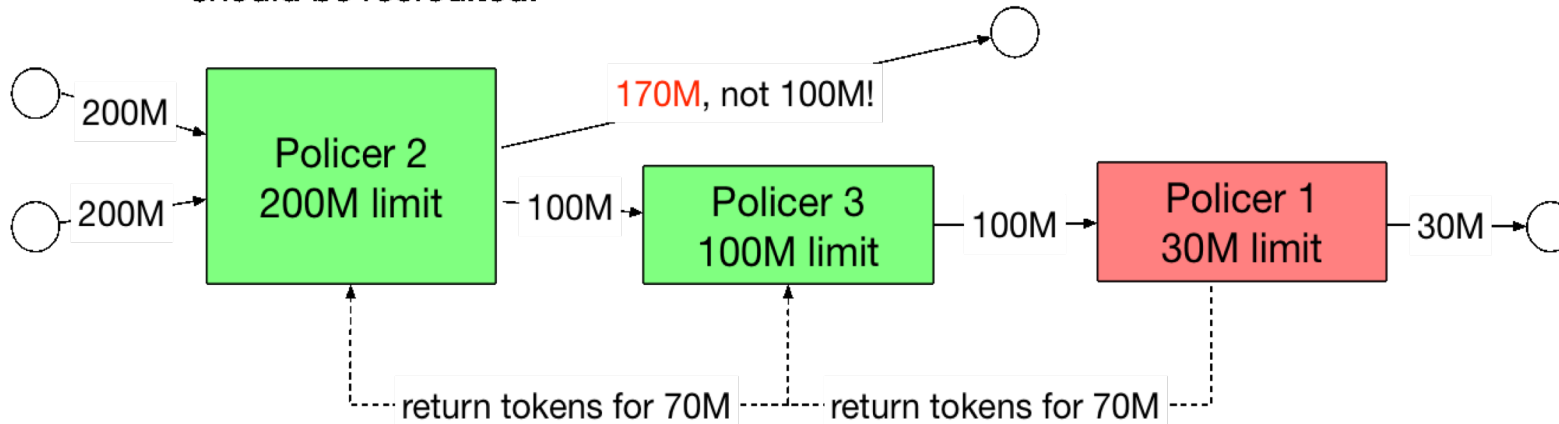
# Traffic policing in eBPF: applying token bucket algorithm

Policers chaining

Scenario 1: The first policer discards a packet, no extra token spent.



Scenario 2: The first two policers has spent tokens on a discarded packet. They should be recredited.



- The output rate must not depend on the order of policers
- If a packet is discarded, recredit the preceding policers
- Policers may not belong to the same logical hierarchy. No common root is required, unlike in qdisc HTB

# Traffic policing in eBPF: applying token bucket algorithm

## Limitations

- Heavy hammer
  - TCP congestion control + token bucket + DROP = capped max rate but underutilized average
  - No buffering: drops are inevitable
- Very thin per TCP flow fairness guarantees
  - No handy TCP session information in tc chain
  - N sub buckets and  $skb \rightarrow hash \% N$
- Token bucket + DSCP remark
  - Only for multi-queue network devices
  - Packets may be received in disorder